



# Neuro-Symbolic Commonsense Knowledge and Reasoning

Yejin Choi

Paul G. Allen School of Computer Science & Engineering  
University of Washington &  
Allen Institute for Artificial Intelligence



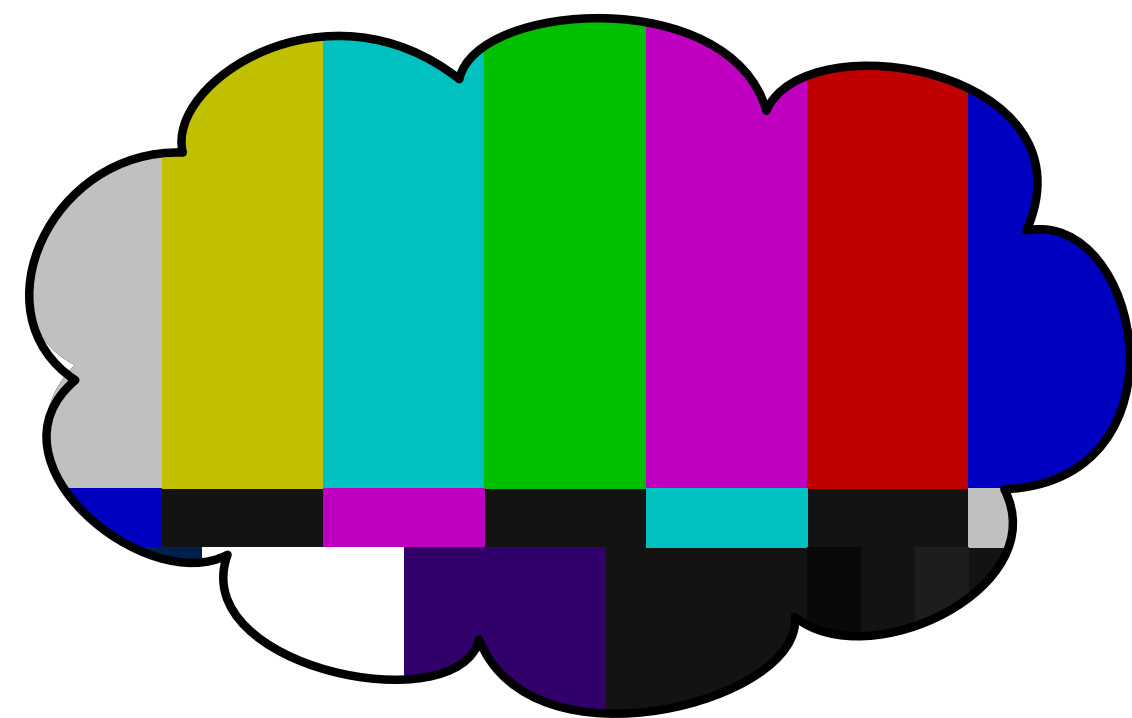
# Despite (super-) human-level performances on leaderboards...

- SOTA neural models are brittle if given adversarial or out-of-domain samples

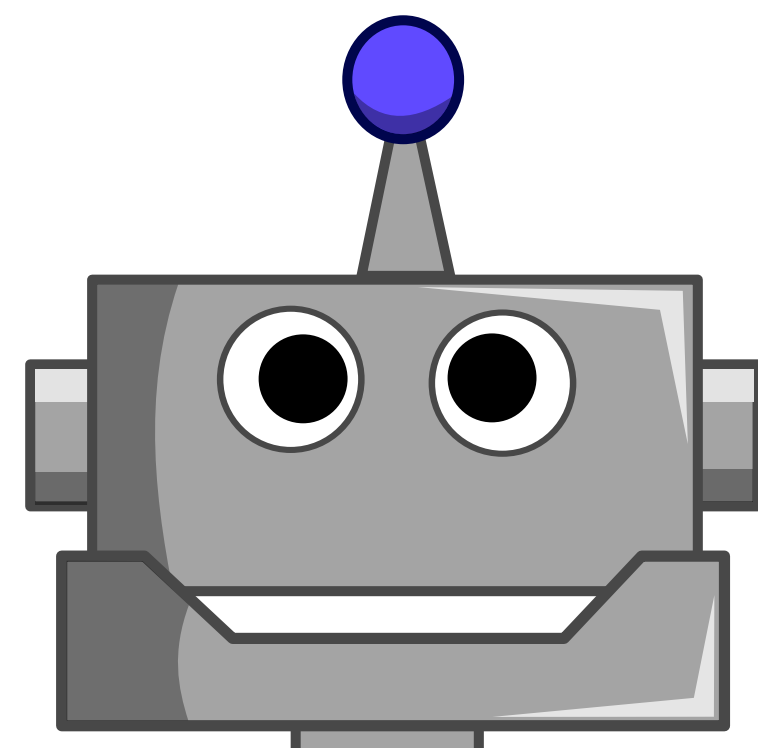
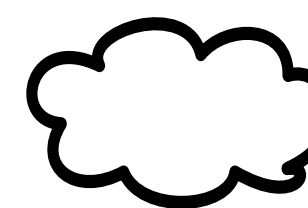
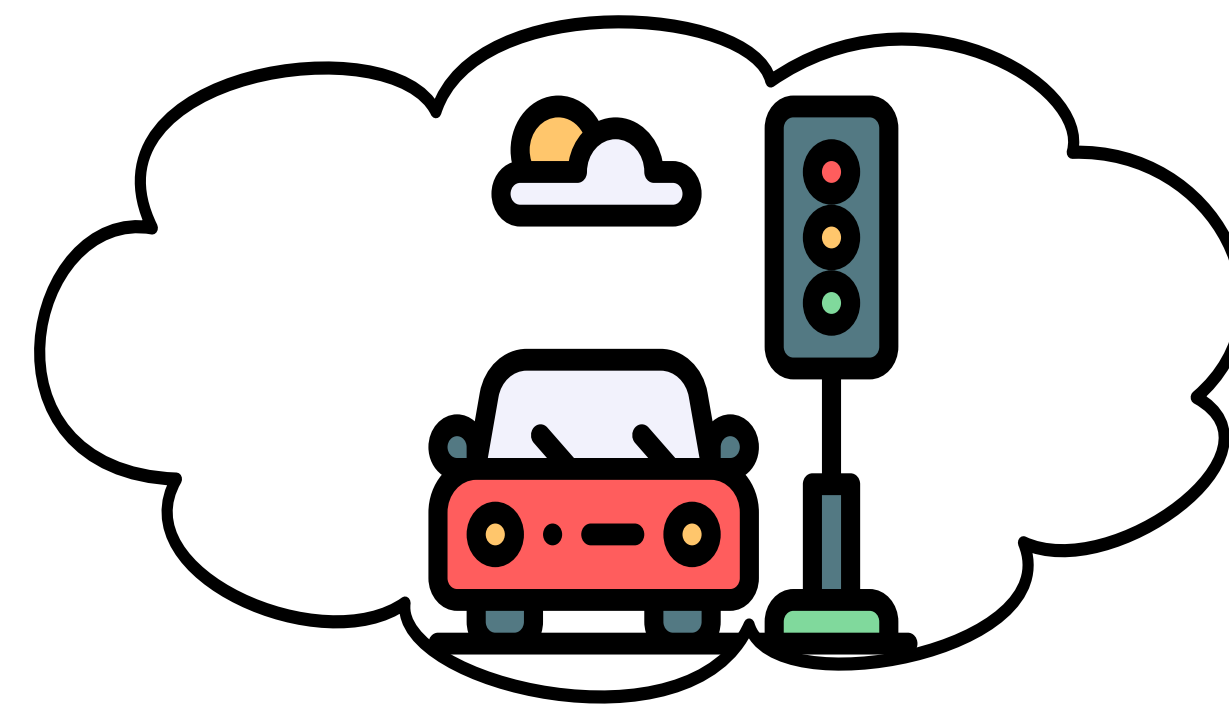
 +  =  Giant panda Object Recognition Szegedy et al, 2014.... Gibbon	<p>.... Nikola Tesla moved to Prague in 1880. ... <b>Tadakatsu moved to Chicago in 1881.</b></p> <p>Where did Tesla move in 1880? <b>Chicago</b></p> <p>QA Jia et al, 2017</p>	 A horse standing in the grass. Captioning MacLeod et al, 2017
---	--	--

Solving only a "dataset"  
without solving the underlying "task"!

(lacking "systematic generalization" (lake & baron 2017, Bahdanau et al & Courville ICLR 2019, Bengio's keynote at NeurIPS 2019))

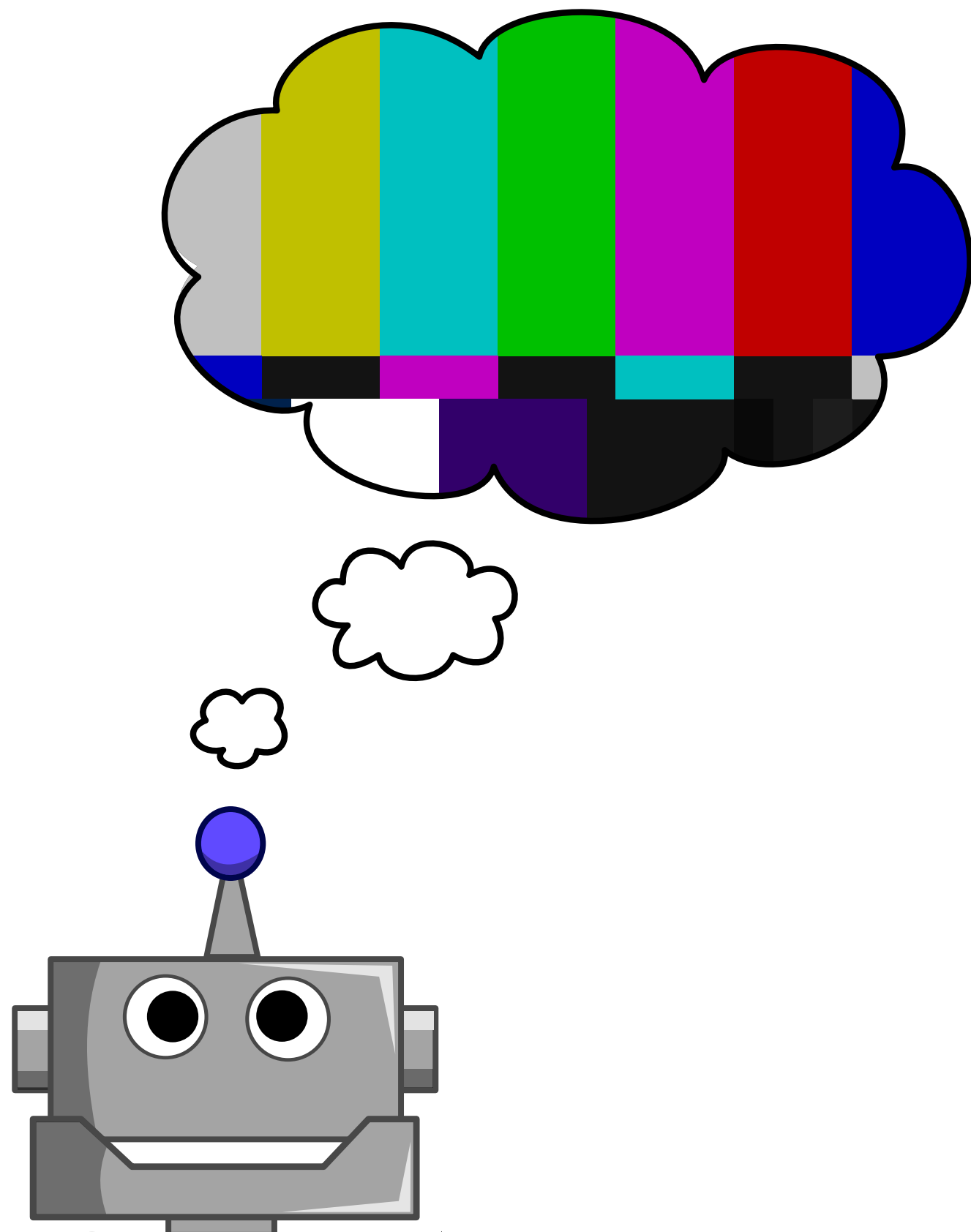


**Let's bridge this gap!**



*Peters et al., 2018;  
Devlin et al., 2018*





Let's bridge this gap!

**SYSTEM 1**  
Intuition & instinct

95%

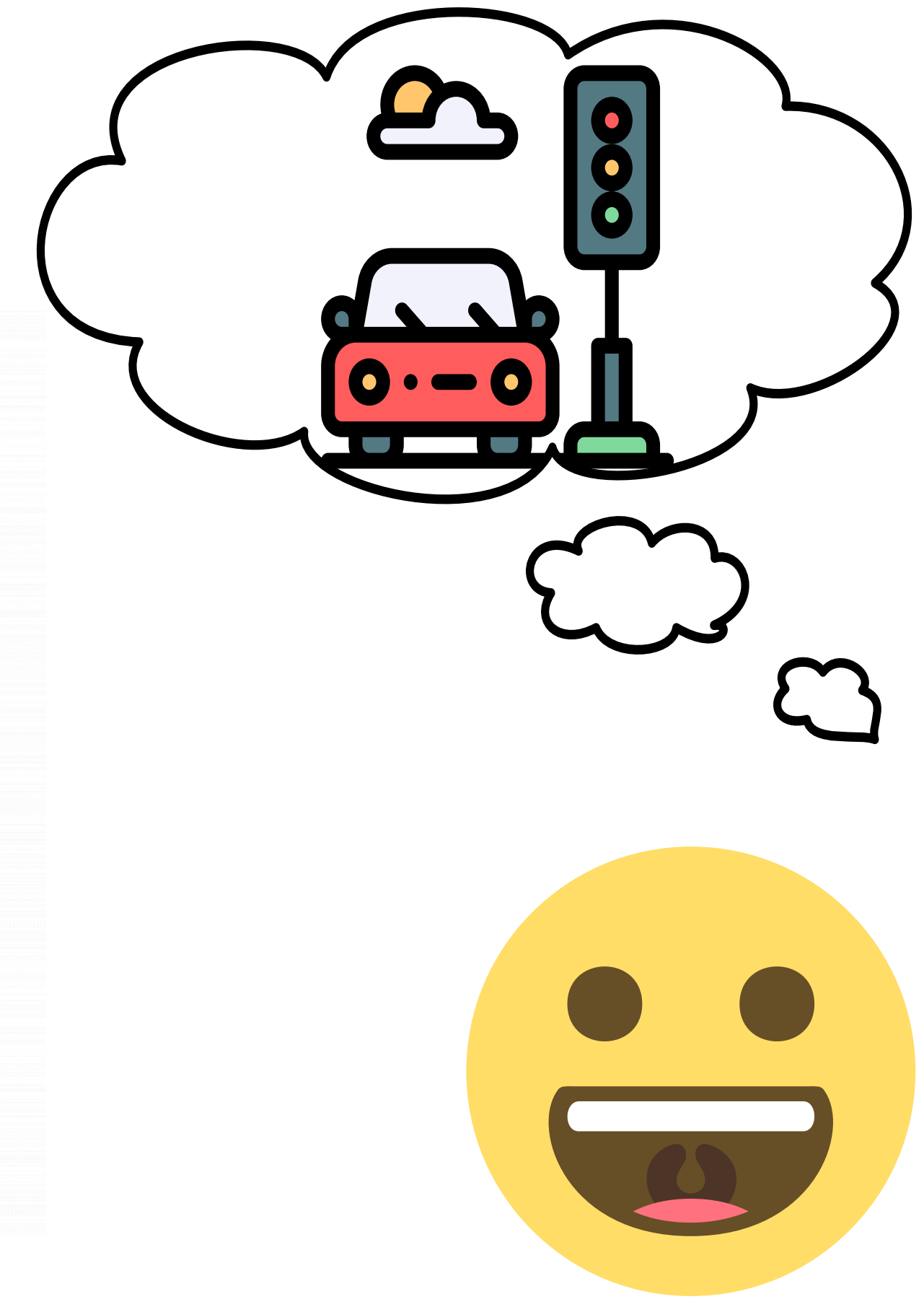
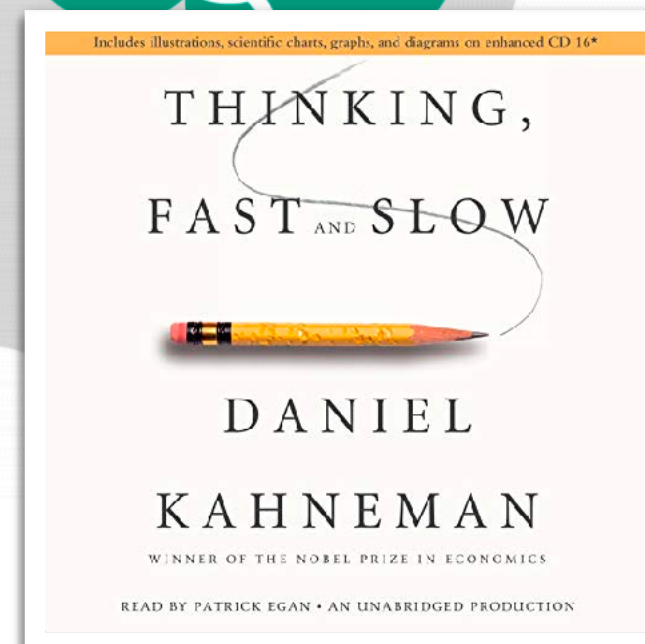
Unconscious  
Fast  
Associative  
Automatic pilot

**SYSTEM 2**  
Rational thinking

5%

Takes effort  
Slow  
Logical  
Lazy  
Indecisive

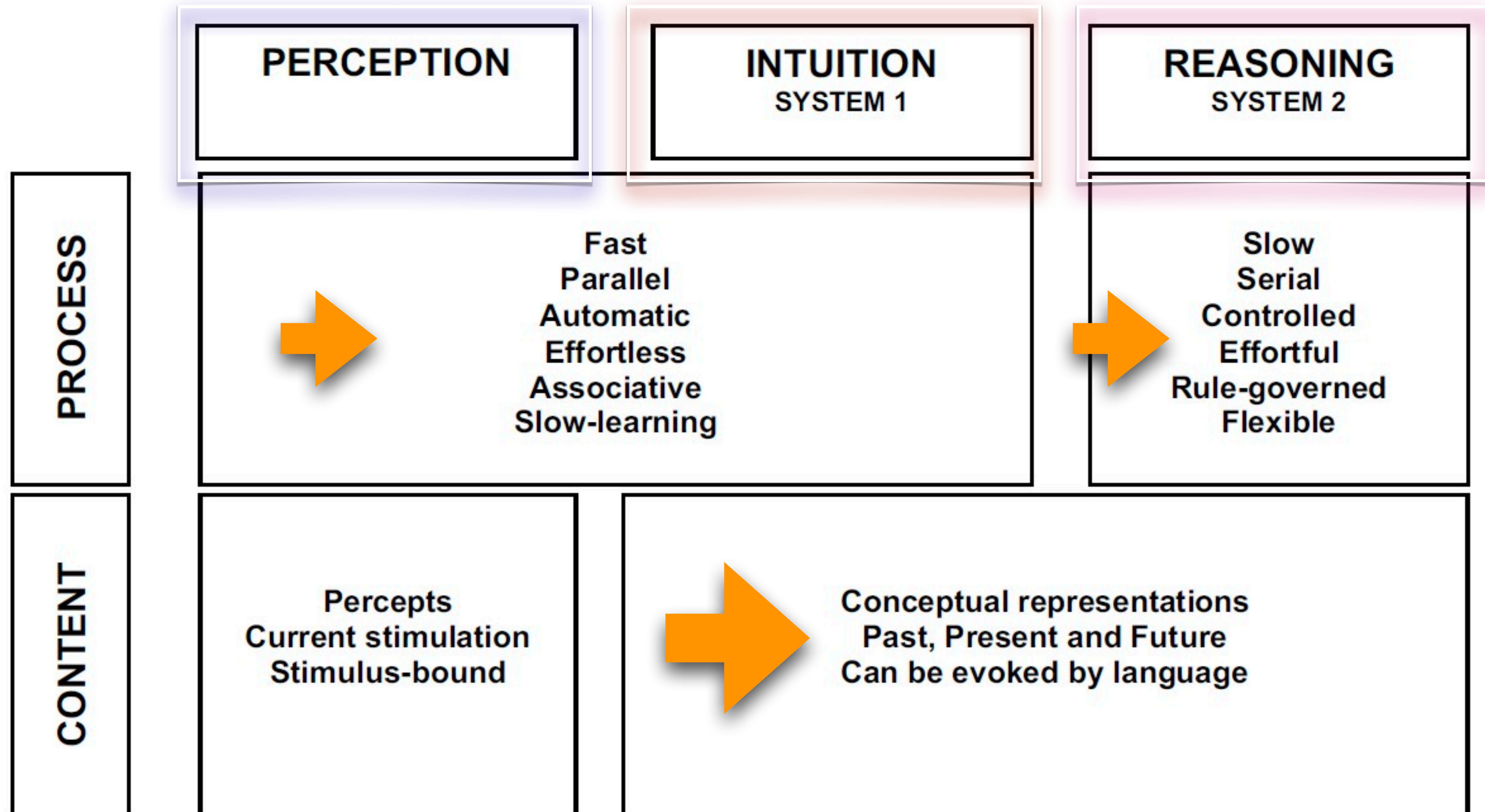
Source: Daniel Kahneman



- **Myth:** we know how to do [system-1 reasoning] with deep learning so we only need to figure out [system-2 reasoning]?

# Kahneman's "three **cognitive** systems"

— "Maps of Bounded Rationality: ..." (Kahneman 2003)





# Kahneman's "three **cognitive** systems"

— "Maps of Bounded Rationality: ..." (Kahneman 2003)

## PERCEPTION

- object recognition
- image segmentation

## INTUITION SYSTEM 1

- Intuitive inferences on
  - pre-conditions and post-conditions
  - what happens before and after?
  - motivations and intents
  - mental and emotional states

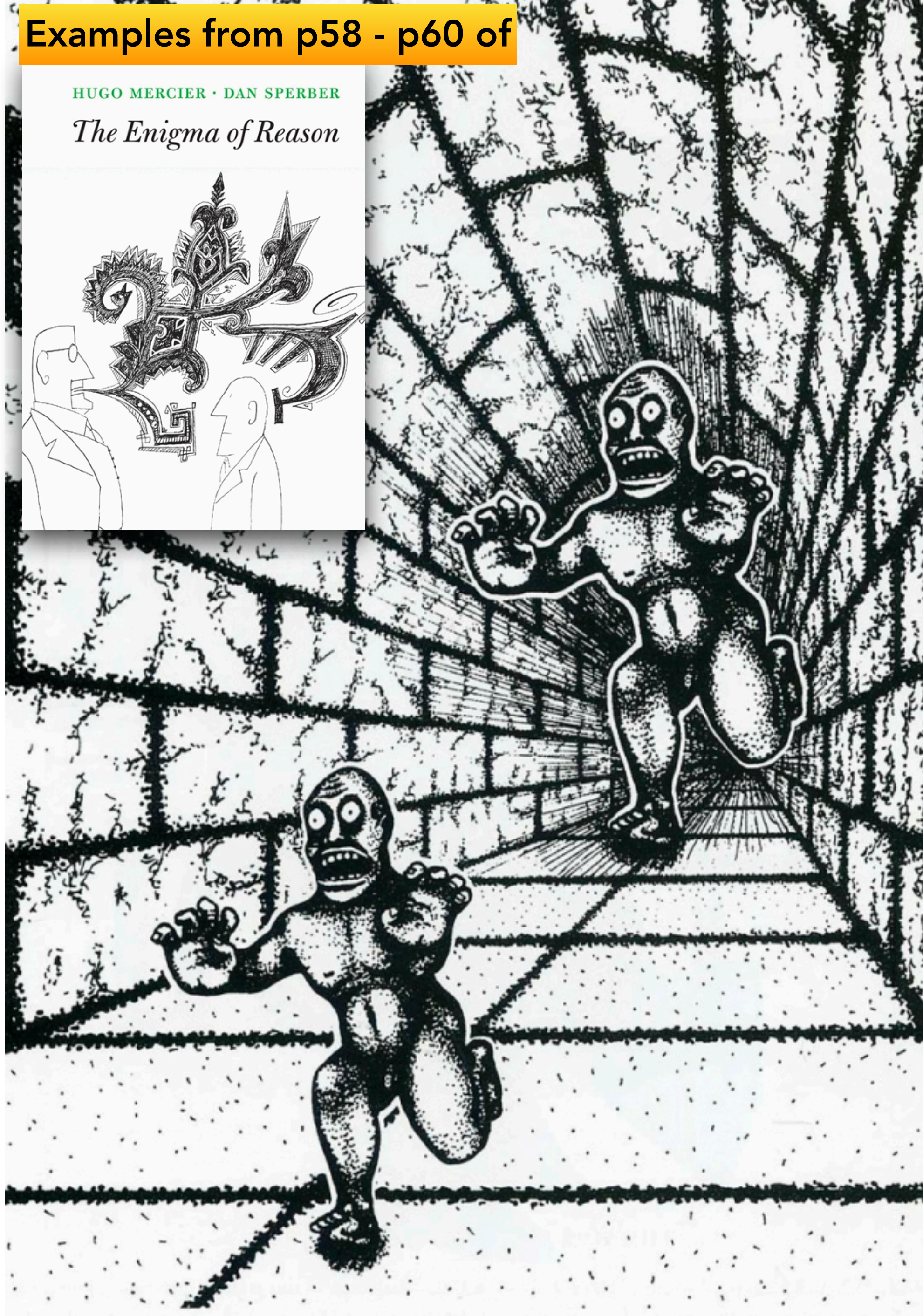
=> This is what humans do every waking minute

## REASONING SYSTEM 2

- solving puzzles
- writing programs
- proving logic theorems
  
- reviewing CVPR papers
- crafting CVPR rebuttals
- giving an invited talk
- writing an op-ed

=> Humans often spend hours (or days) **not** doing this sort of reasoning at all...





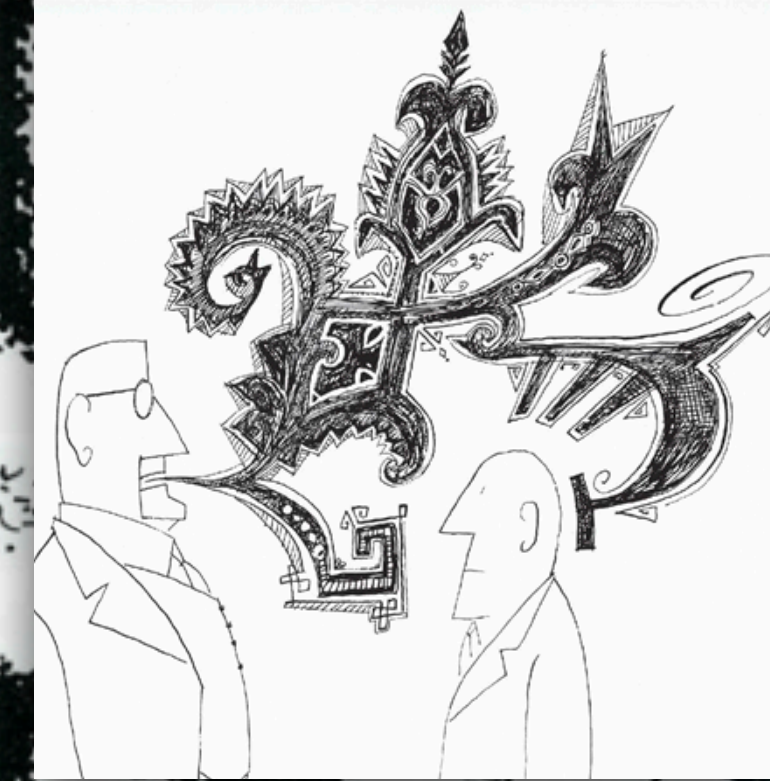
# Roger Shepard's "monsters in a tunnel"

- **Two monsters are running** (rather than standing still on one foot)
- **One is chasing another** (rather than trying to copy his movements)
- **The chaser has hostile intentions and the chased is afraid** (even though two faces are identical)

INTUITION  
SYSTEM 1

- Intuitive inferences on
  - pre-conditions and post-conditions
  - what happens before and after?
  - motivations and intents
  - mental and emotional states





# Roger Shepard's "monsters in a tunnel"

- **Two monsters are running** (rather than standing still on one foot)
- **One is chasing another** (rather than trying to copy his movements)
- **The chaser has hostile intentions and the chased is afraid** (even though two faces are identical)

## INTUITION SYSTEM 1

### Important Observations:

- None of these inferences is absolutely true. The inferences are **stochastic** in nature. Everything is **defeasible** with additional context.
- A great deal of **intuitive inferences** are **commonsense inferences**, a great deal of which can be best described in **natural language** — **full scope of language**, **not just words**, or even **graphs of words**

- **Intuitive inferences on**
  - **pre-conditions and post-conditions**
  - **what happens before and after?**
  - **motivations and intents**
  - **mental and emotional states**



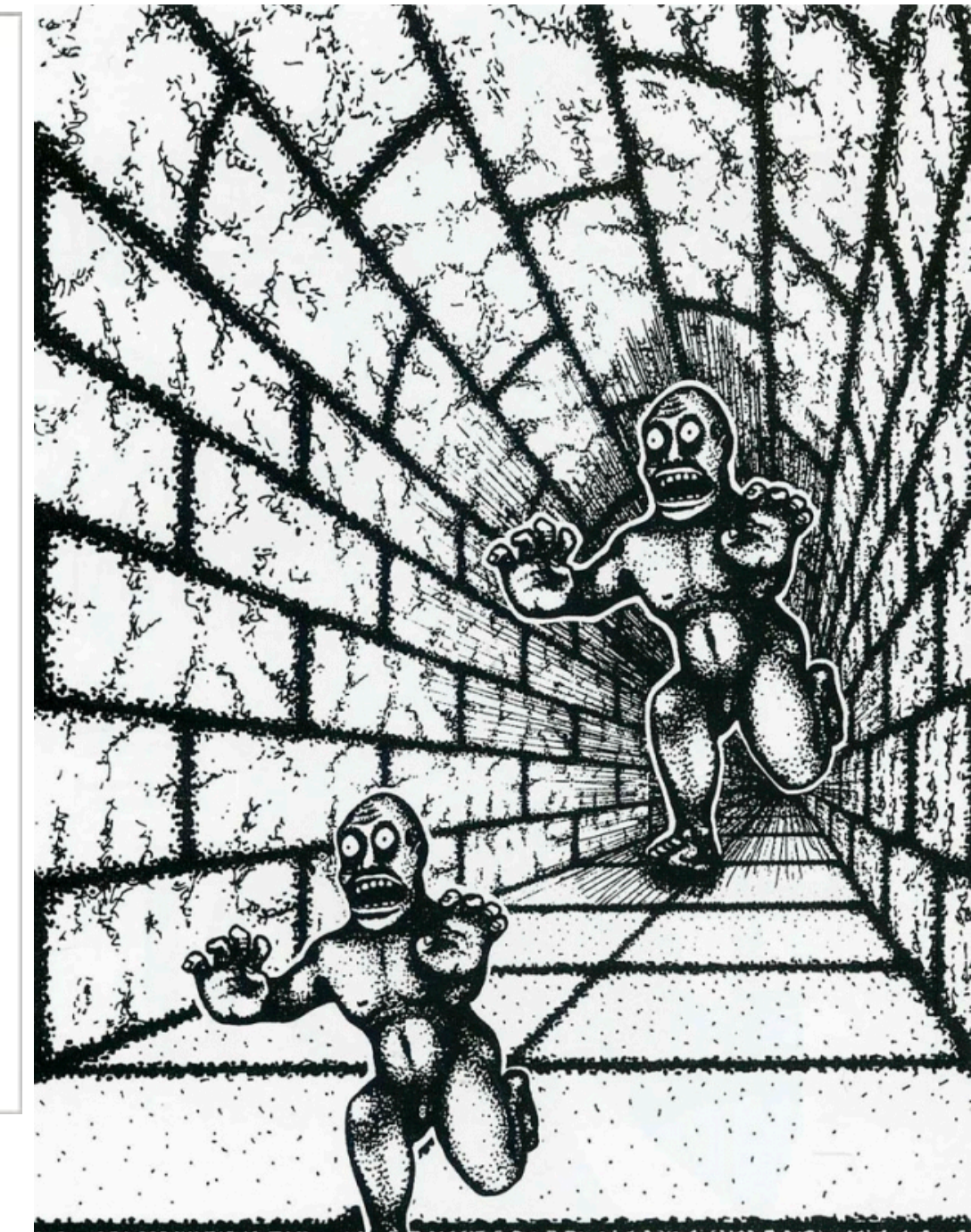
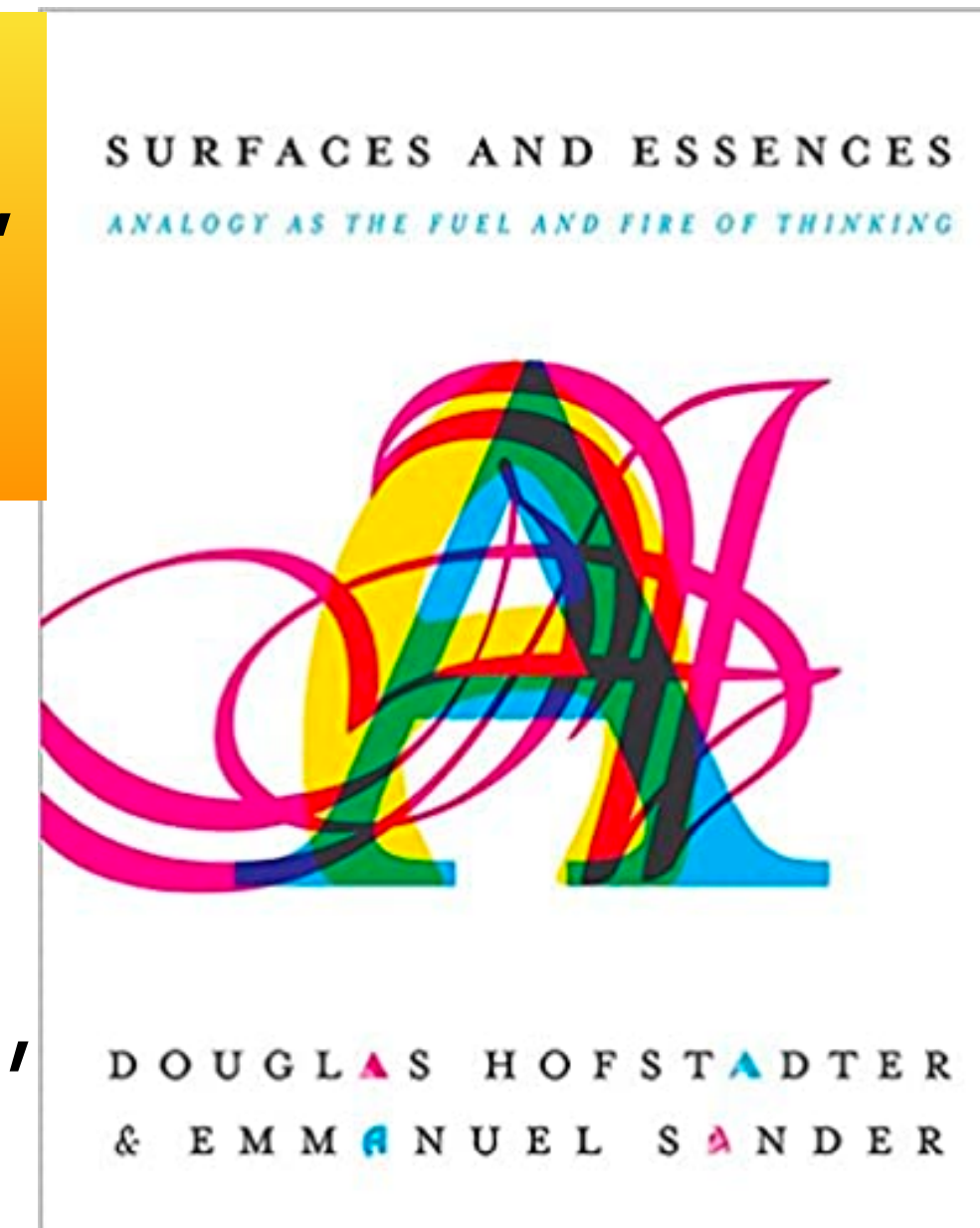
# Tldr; 🔥 Language and Symbols 🔥

## 🔥 Reasoning as Generation 🔥

### 1. *Language* as the *symbols*

- all of it —
- not just **words** (ImageNet labels, ...)
- not just small **sets of words** (scene graphs,

"Categories (concepts) vastly outnumber words, and require free-form open text descriptions"



### 2. Reasoning as *generative tasks*

- As opposed to *discriminative tasks* (i.e., categorization)
- Because the space of reasoning in language is **infinite**

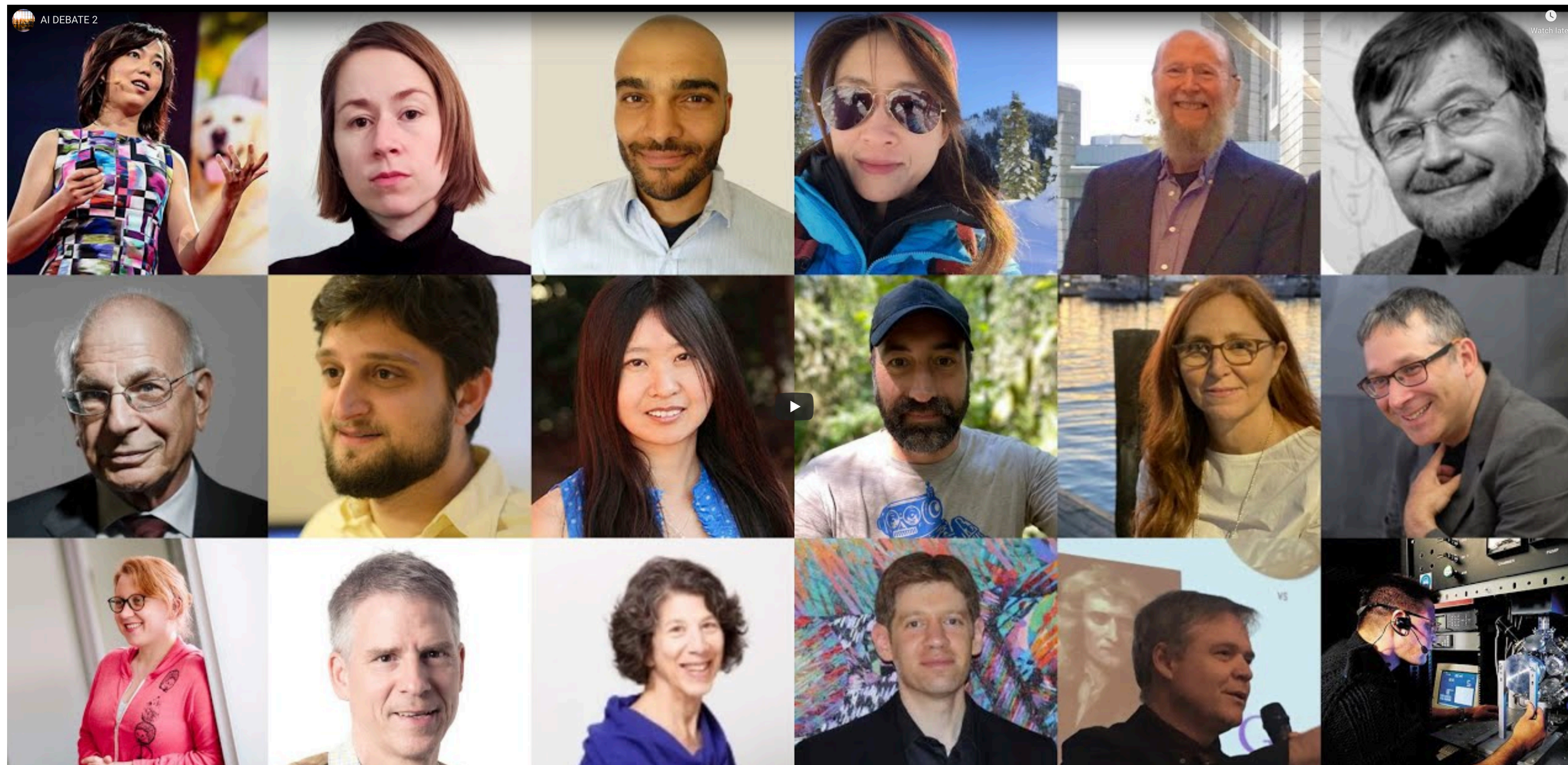
"thinking out loud"

We often think as we speak,  
on the fly, word-by-word  
without enumerating all possible  
alternative sentences



# AI Debate 2 at montreal.ai

- <https://montrealartificialintelligence.com/aidebate2.html>
  - **Daniel Kahneman:** <https://youtu.be/2zNd69ZGZ8o>
- <https://venturebeat.com/2021/01/02/leading-computer-scientists-debate-the-next-steps-for-ai-in-2021/>
- Stay tuned for AI Debate 3 (Dec 23 2021!)





# In this talk: Reasoning as Generation

- **Part 1:** unsupervised inference-time algorithms

Reasoning thru  
**Neural Backpropagation**


DeLorean

Reasoning thru  
**Search with Logical Constraints**

NeuroLogic

Reasoning thru  
**Distributional Neural Imagination**

Reflective Decoding

=> How to make  from (off-the-shelf) neural language models

- **Part 2:** supervised, but with declarative knowledge

?

?

?

- **Part 3:** benchmarks and algorithmic bias reduction





Back to the Future:

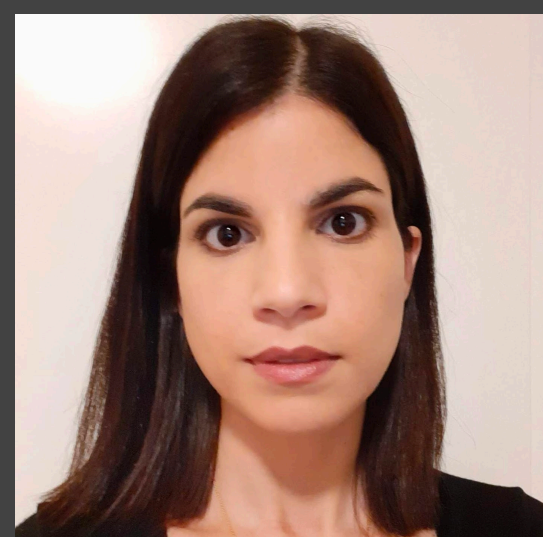
# Unsupervised Backprop-based Decoding for Counterfactual and Abductive Commonsense Reasoning

EMNLP 2020

Lianhui Qin



Vered  
Shwartz



Peter  
West



Chandra  
Bhagavatula



Jena  
Hwang



Ronan  
LeBras



Antoine  
Bosselut



Me



# Abductive Reasoning

(Bhagavatula et al., 2019)



## *Past Observation*

Ray hung a tire on a rope to make his daughter a swing.

# What happened in between?



## *Future Observation*

Ray ran to his daughter to make sure she was okay.



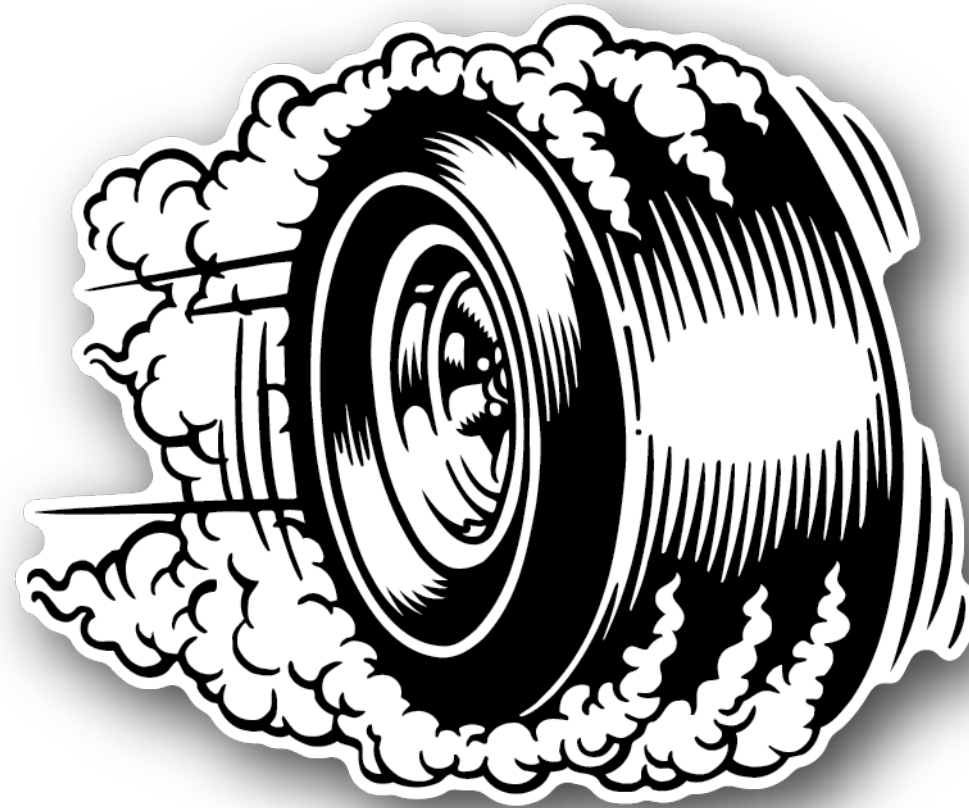
# Abductive Reasoning

(Bhagavatula et al., 2019)



*Past Observation*

Ray hung a tire on a rope to make his daughter a swing.



*Hypothesis*

*She hit the rope and the tire fell on top of her.*



*Future Observation*

Ray ran to his daughter to make sure she was okay.

**Abductive Reasoning** = Inference to the best explanation to partial observation (Peirce 1960)  
**Abduction** != Induction or Deduction

# Abductive Reasoning

(Bhagavatula et al., 2019)

## *Past Observation*

Ray hung a tire on a rope to make his daughter a swing.

## *Hypothesis*

*She hit the rope and the tire fell on top of her.*

## *Future Observation*

Ray ran to his daughter to make sure she was okay.

# Counterfactual Reasoning

(Qin et al., 2019)



An example from the  
“TimeTravel” dataset  
(Qin et al., EMNLP 2019)

Zeke was throwing a party.

All his friends were dressing up for this Halloween party. Story context changes...

Zeke thought about being a vampire or a wizard.

Then he decided on a scarier costume.

Zeke dressed up like a skeleton.

What if this is a **Game of Thrones themed party** instead of a **Halloween party**? 🤔







An example from the  
“TimeTravel” dataset  
(Qin et al., EMNLP 2019)

Story ending doesn't  
make sense now...

Zeke was throwing a party.

All his friends were dressing up for this Halloween party.

Story context changes...

Zeke thought about being a vampire or a wizard.

Then he decided on a scarier costume.

Zeke dressed up like a skeleton.

What if this is a **Game of Thrones themed party** instead of a **Halloween party**?





An example from the  
“TimeTravel” dataset  
(Qin et al., EMNLP 2019)

Story ending doesn't  
make sense now...

Zeke was throwing a party.

All his friends were dressing up for this Halloween party.

Story context changes...

Zeke thought about being a vampire or a wizard.

Then he decided on a scarier costume.

Zeke dressed up like a skeleton.

What if this is a **Game of Thrones themed party** instead of a **Halloween party**?



## Counterfactual Reasoning

(Qin et al., 2019)

Reasoning about the alternative **future**  
based on counterfactual **past**.



An example from the  
“TimeTravel” dataset  
(Qin et al., EMNLP 2019)

Story ending doesn't  
make sense now...

Zeke was throwing a party.

All his friends were dressing up for this Halloween party.

Story context changes...

Zeke thought about ~~being a vampire or a wizard.~~

~~Then he decided on a scarier costume.~~

Zeke dressed up like a ~~skeleton.~~

What if this is a **Game of Thrones themed party** instead of a **Halloween party**?



Zeke was throwing a party.

All his friends were dressing up for this Halloween party.

Only do minimal edit!

Zeke thought about *Lannister, but he didn't want to look like a Lannister.*

Consistent  
now!!

*He wanted to look like a Stark.*

Zeke dressed up like a *Stark.*



# Abductive Reasoning

(Bhagavatula et al., 2019)

## Past Observation

Ray hung a tire on a rope to make his daughter a swing.

## Hypothesis

*She hit the rope and the tire fell on top of her.*

## Future Observation

Ray ran to his daughter to make sure she was okay.

## Story Context

Zeke was throwing a party.

All his friends were dressing up for this Halloween party.

All his friends were dressing up for this Game of Thrones themed party.

## Rewritten Ending

*Zeke thought about Lannister, but he didn't want to look like a Lannister.*

*He wanted to look like a Stark.*

*Zeke dressed up like a Stark.*

## Original Ending

Zeke thought about being a vampire or a wizard.

Then he decided on a scarier costume.

Zeke dressed up like a skeleton.

# Counterfactual Reasoning

(Qin et al., 2019)

# Abductive Reasoning

(Bhagavatula et al., 2019)

Both involve ***nonmonotonic reasoning*** with  
past context  $X$  and future constraint  $Z$



# Counterfactual Reasoning

(Qin et al., 2019)



Pretrained Language Models are successful on many tasks...

How are Pretrained LMs on  
the Nonmonotonic Reasoning?

Let's first see the abductive case...

*Y*

The little girl liked it and was thrilled at it.

Pre-trained GPT2

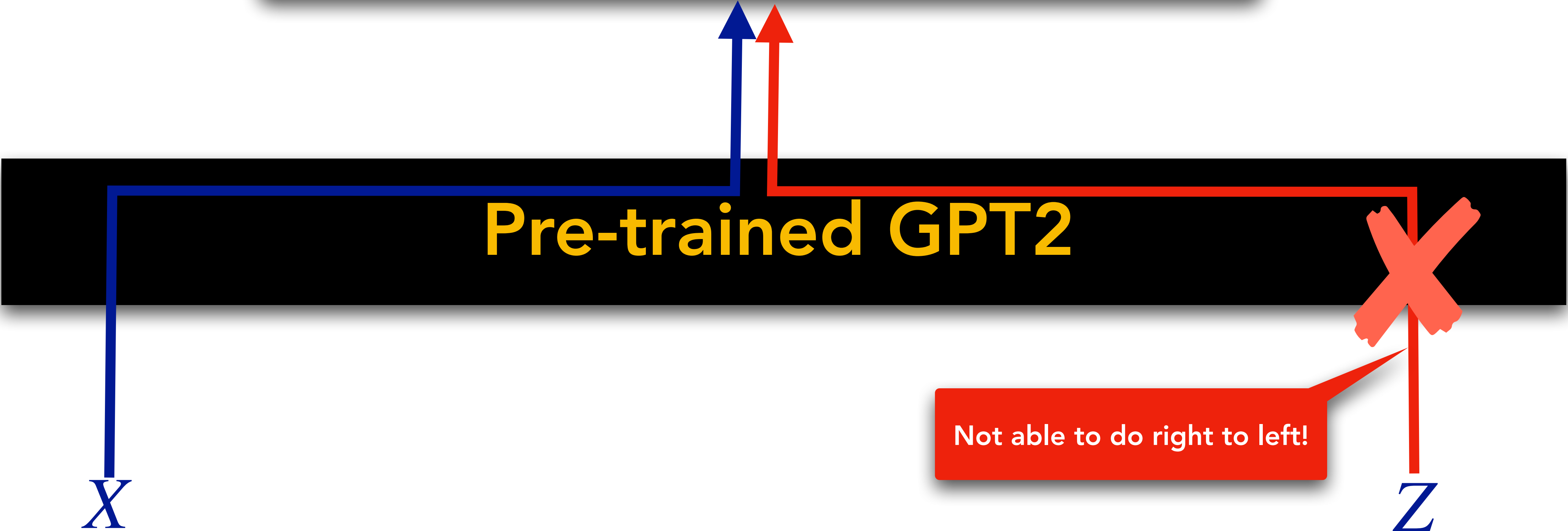
*X*

Ray hung a tire on a rope to make his daughter a swing.

*Z*

Ray ran to his daughter to make sure she was okay.

Not able to do right to left!



Why not just concatenate both direction?

$Y$

But the swing didn't go off, so they moved down the slope towards the



Doesn't make sense!

Pre-trained GPT2

$Z[s]X$

Ray ran to his daughter to make sure she was okay.

Ray hung a tire on a rope to make his daughter a swing.



Try again?

$Y$

As the swing moved, the girl's cries sounded in his ears.



Doesn't make sense!

Pre-trained GPT2

$Z[s]X$

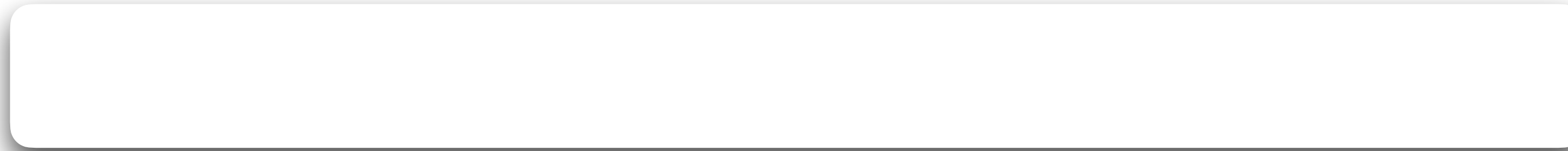
Ray ran to his daughter to make sure she was okay.

Ray hung a tire on a rope to make his daughter a swing.

Something might have been missing here...



$Y$



Backpropagation!!

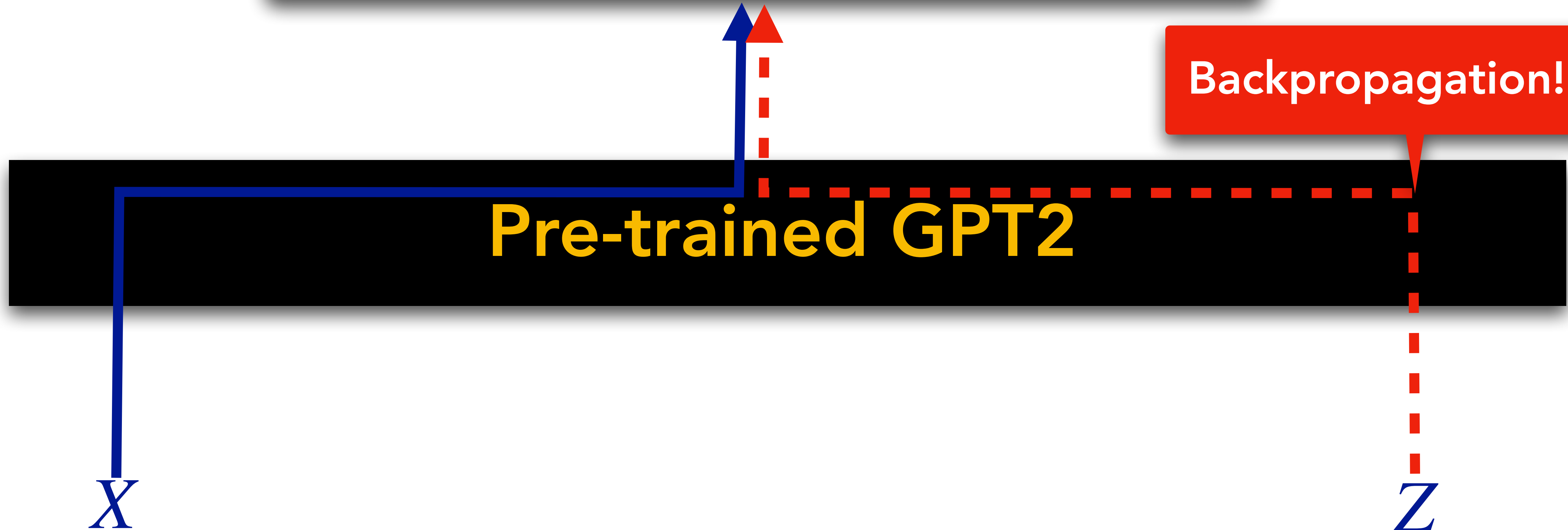
Pre-trained GPT2

$X$

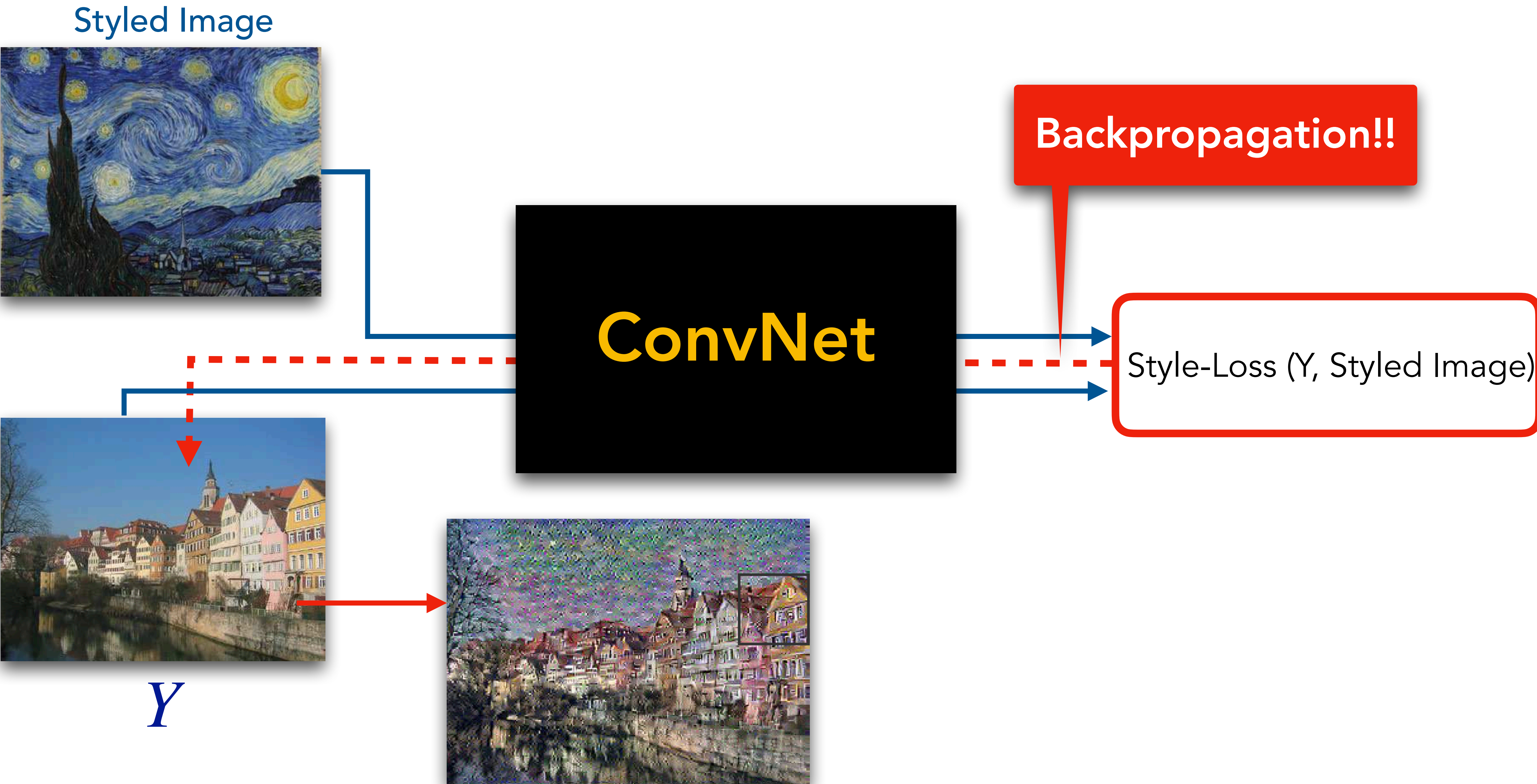
Ray hung a tire on a rope to make his daughter a swing.

$Z$

Ray ran to his daughter to make sure she was okay.

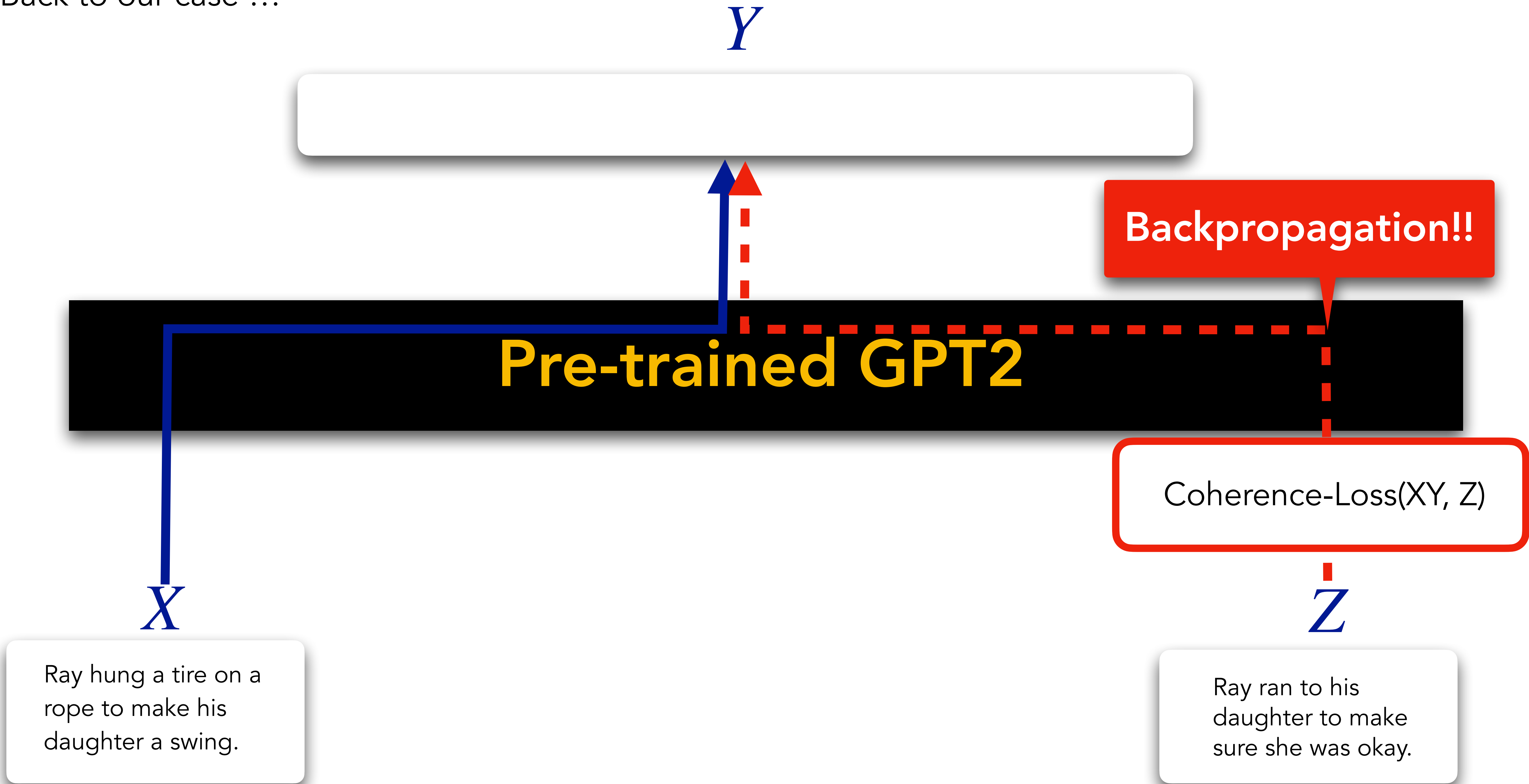


Inspired by “Image Style Transfer” (Gatys et al, 2016)...





Back to our case ...





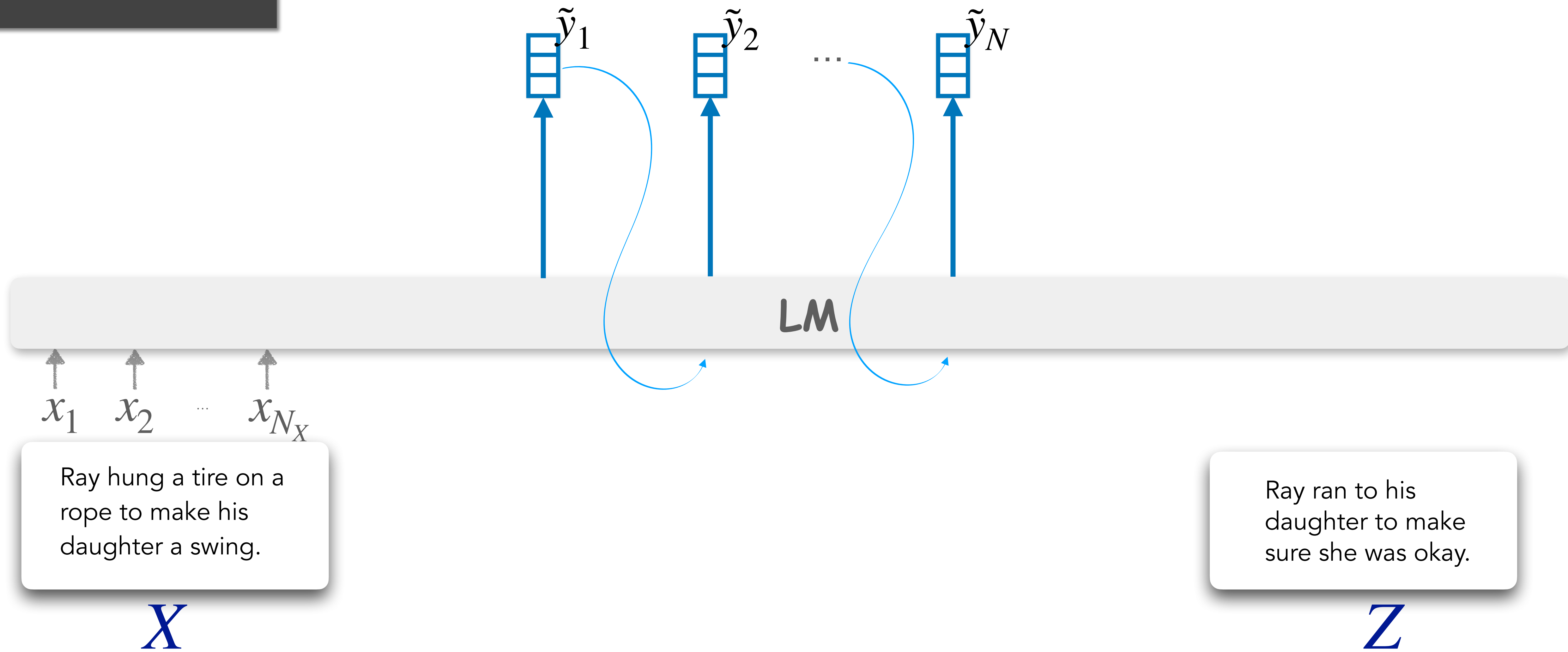
# DELOREAN

(DEcoding for nonmonotonic LOgical REAsoNing)



## Initialization

Just as how you do  
regular decoding



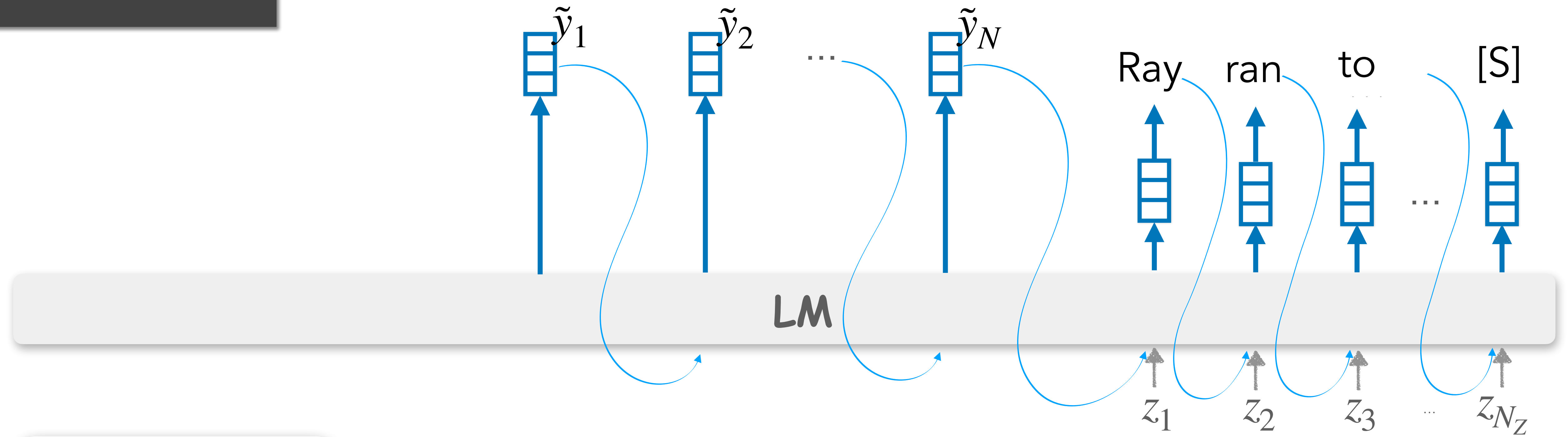




# Backward Pass

Backpropagate  
future information

$Y$



Ray hung a tire on a  
rope to make his  
daughter a swing.

$X$

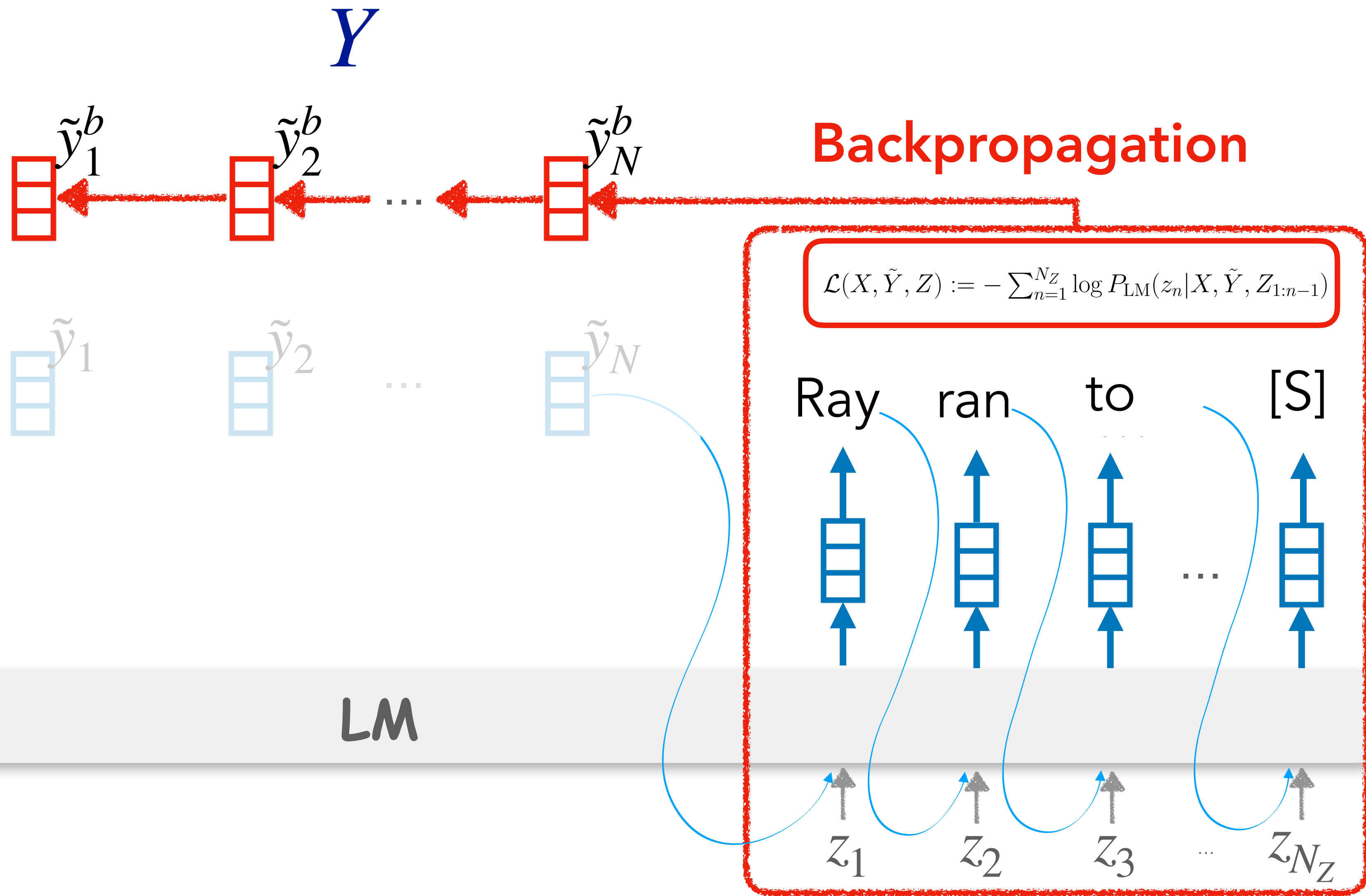
Ray ran to his  
daughter to make  
sure she was okay.

$Z$



# Backward Pass

Backpropagate  
future information  
 $\text{Loss}(Z | X, Y)$



Ray hung a tire on a  
rope to make his  
daughter a swing.

$X$

Ray ran to his  
daughter to make  
sure she was okay.

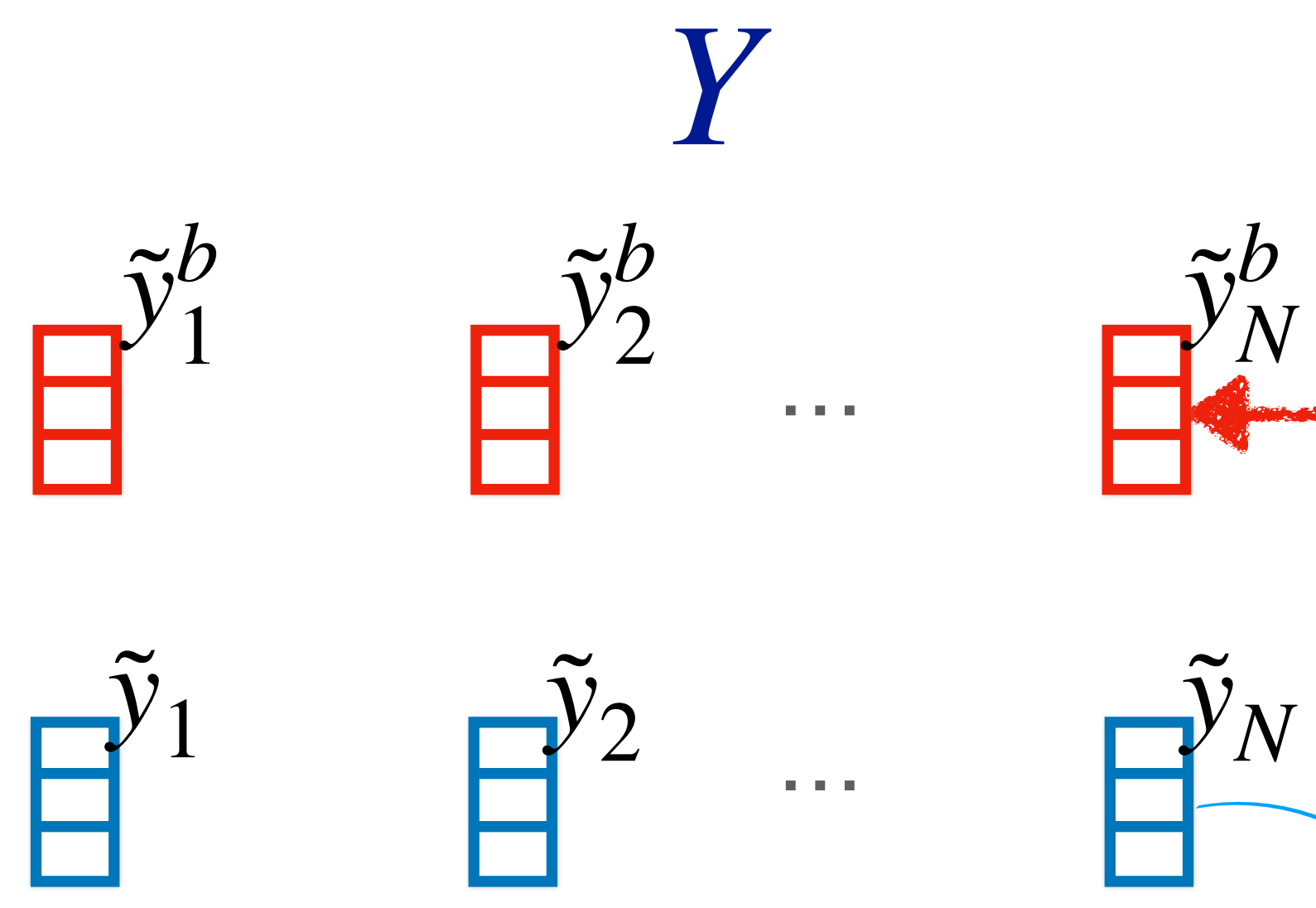
$Z$





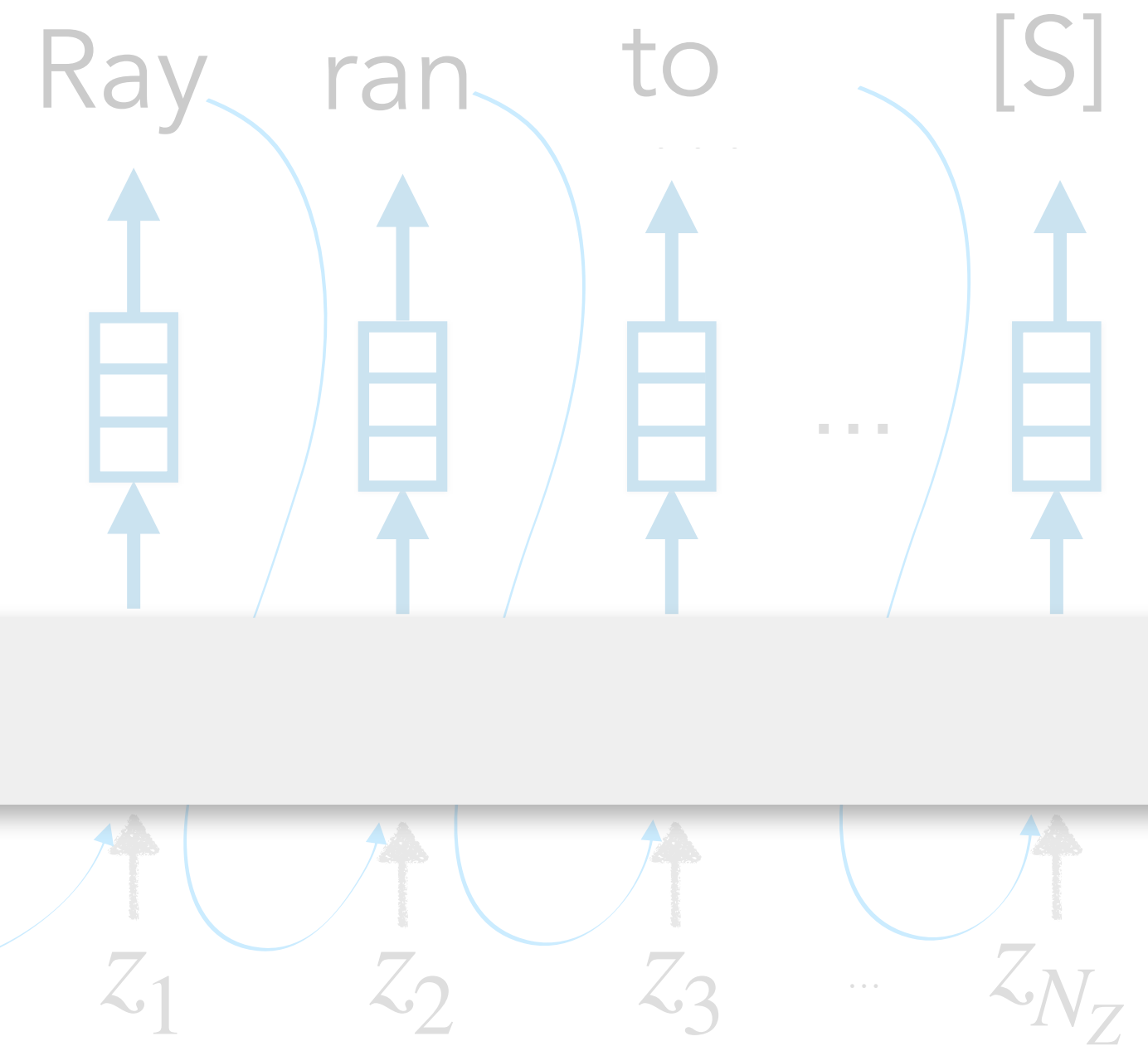
# Forward Pass

Mix both past and future information



# Backpropagation

$$\mathcal{L}(X, \tilde{Y}, Z) := - \sum_{n=1}^{N_Z} \log P_{LM}(z_n | X, \tilde{Y}, Z_{1:n-1})$$



Ray hung a tire on a rope to make his daughter a swing.

Ray ran to his daughter to make sure she was okay.

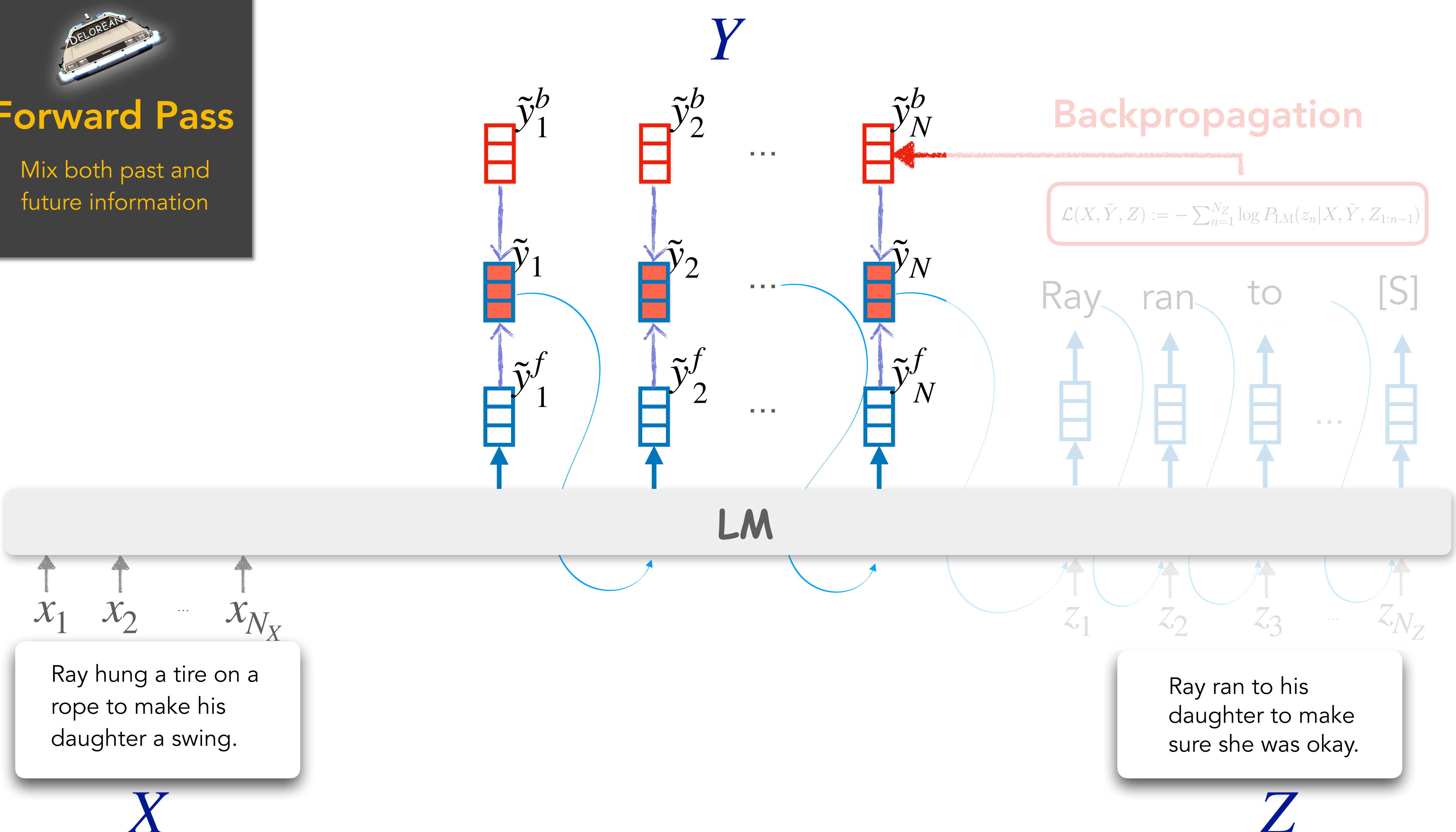
$X$

$Z$



# Forward Pass

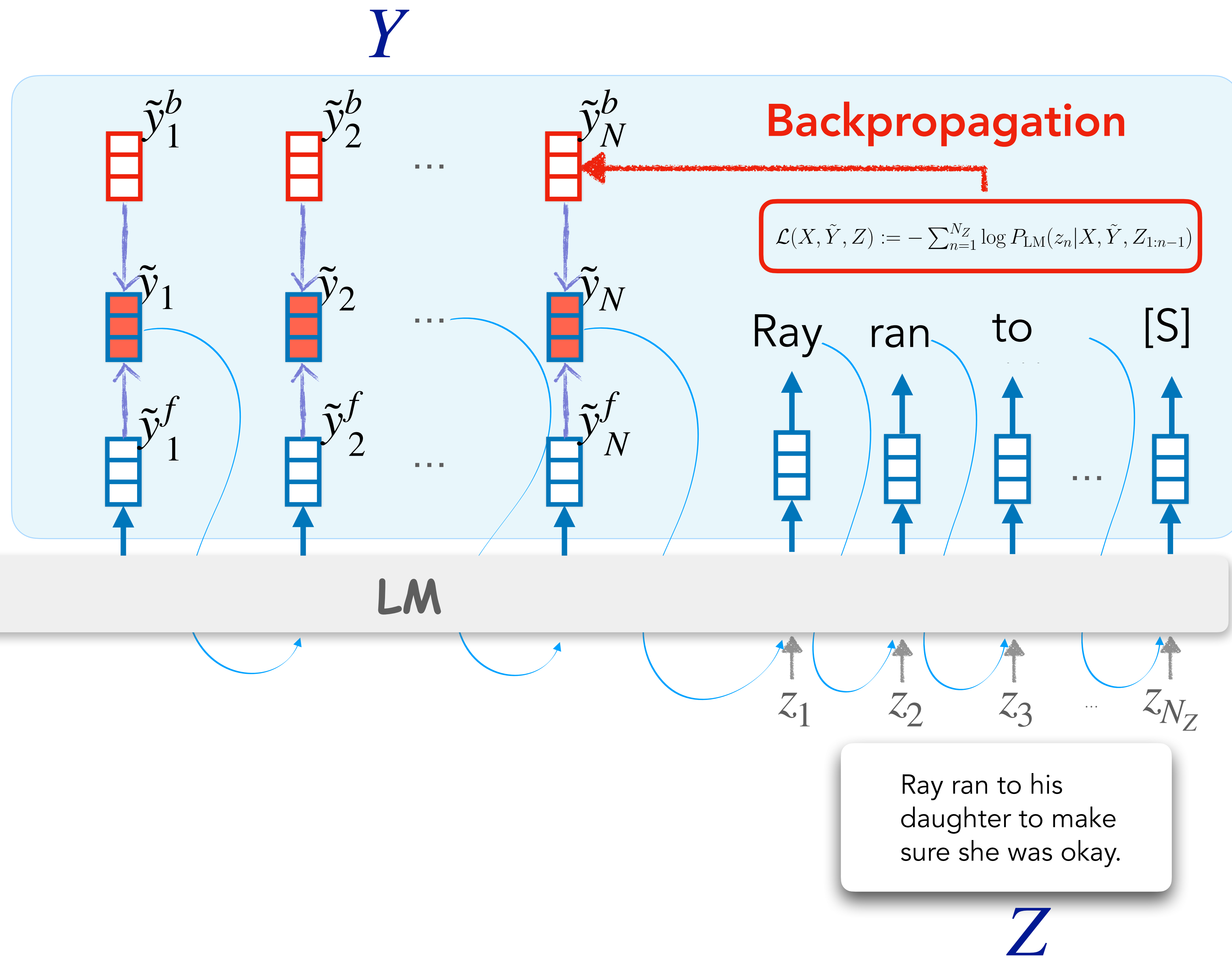
Mix both past and future information







Repeat  
 $T$  times





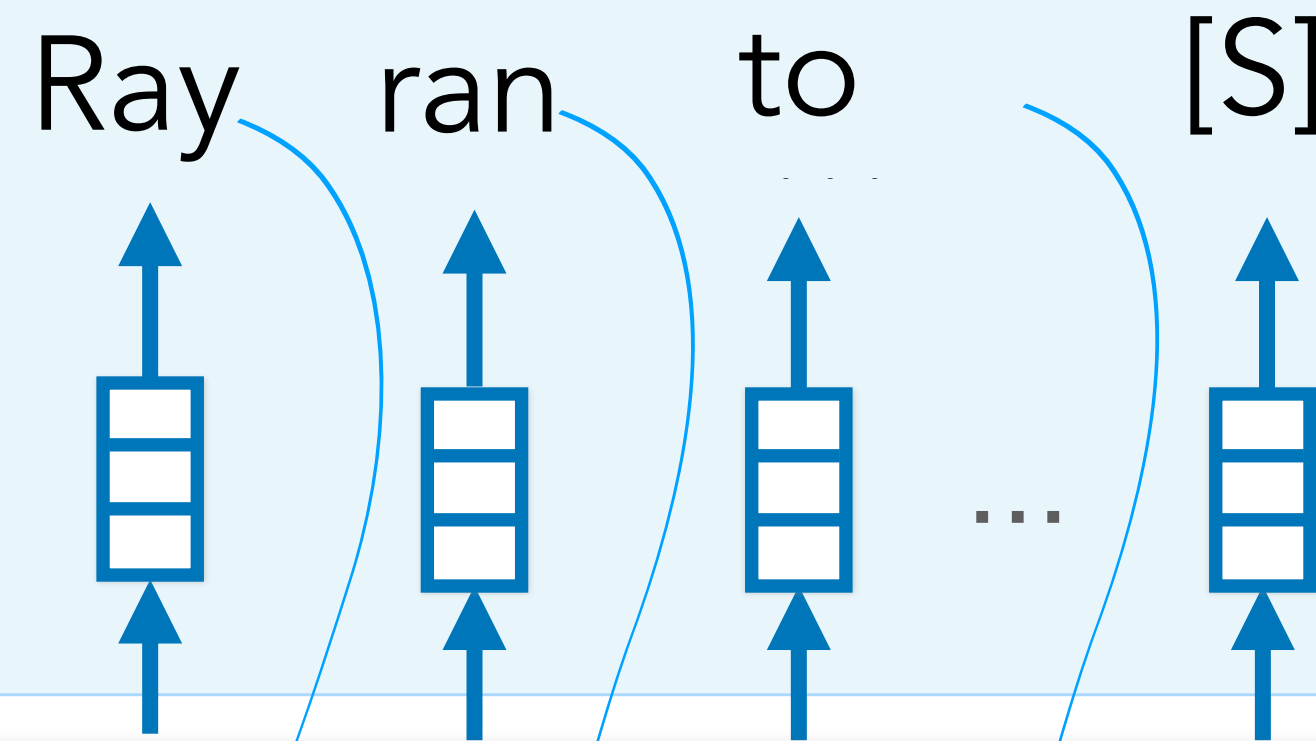
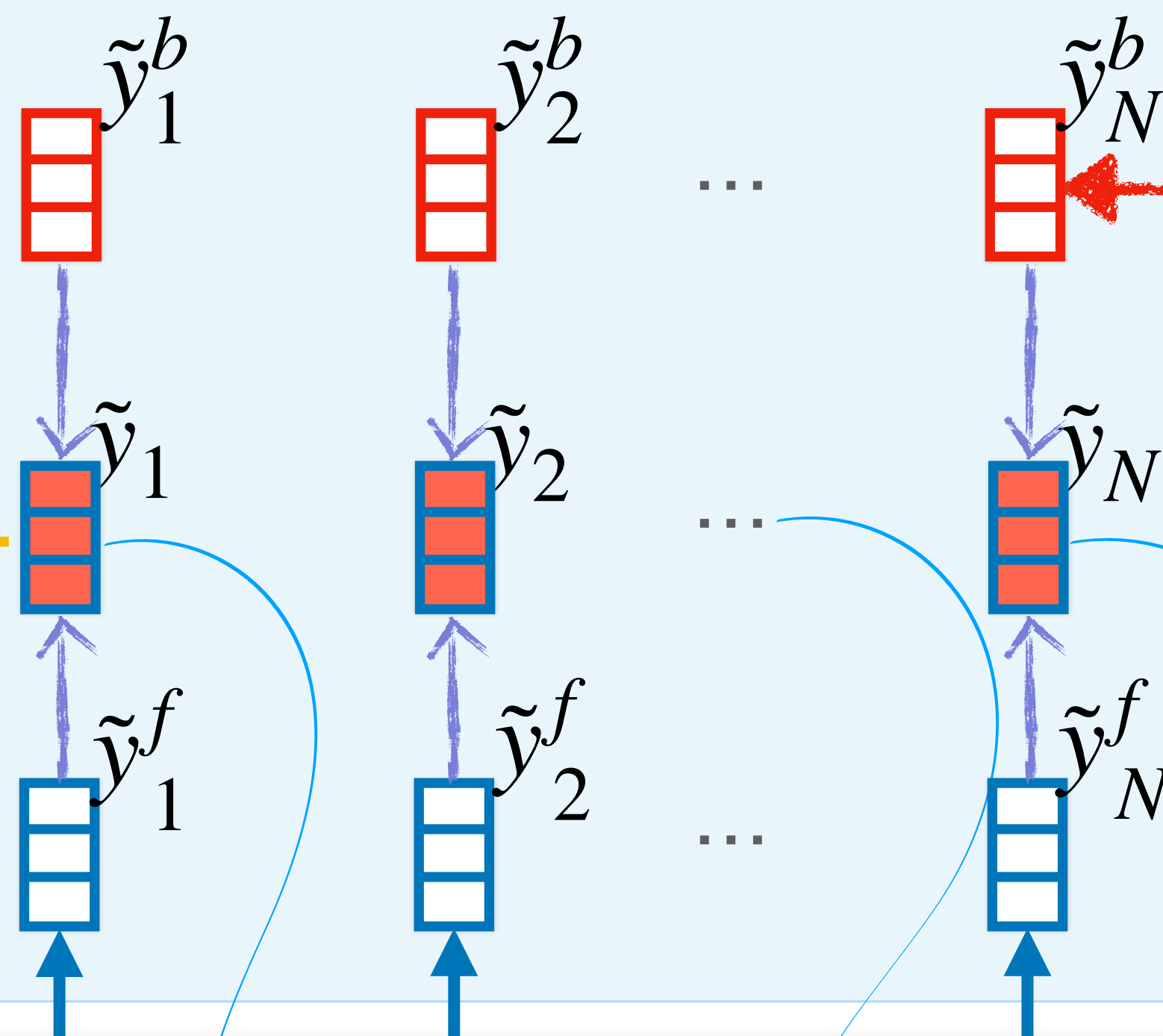
Sampling

$Y$

Output: She hit the rope and the tire fell on top of her.

Backpropagation

$$\mathcal{L}(X, \tilde{Y}, Z) := - \sum_{n=1}^{N_Z} \log P_{\text{LM}}(z_n | X, \tilde{Y}, Z_{1:n-1})$$



$x_1$   $x_2$  ...  $x_{N_X}$

Ray hung a tire on a rope to make his daughter a swing.

$X$

$z_1$   $z_2$   $z_3$  ...  $z_{N_Z}$

Ray ran to his daughter to make sure she was okay.

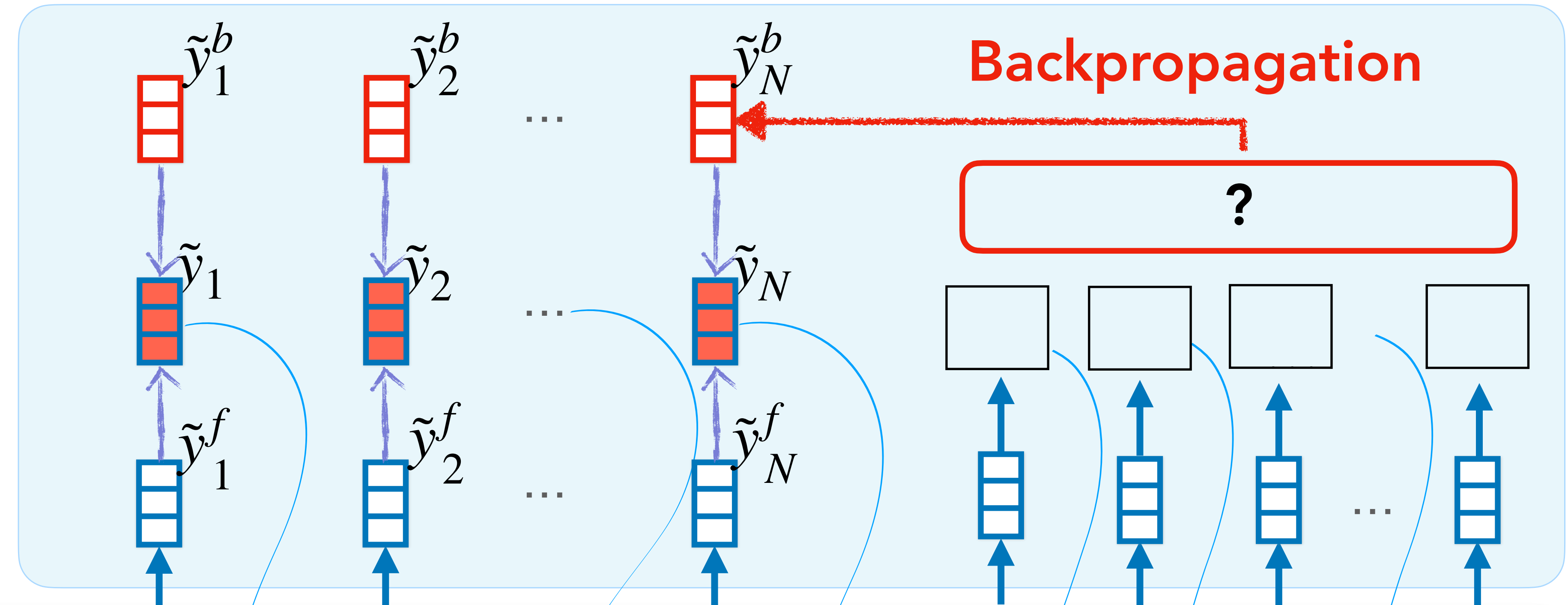
$Z$





# Counterfactual Reasoning?

$Y$



$\tilde{y}_1^b$   $\tilde{y}_2^b$  ...  $\tilde{y}_N^b$

$\tilde{y}_1^f$   $\tilde{y}_2^f$  ...  $\tilde{y}_N^f$

$\tilde{y}_1^f$   $\tilde{y}_2^f$  ...  $\tilde{y}_N^f$

Backpropagation

?

LM

$x_1$   $x_2$  ...  $x_{N_X}$

$z_1$   $z_2$   $z_3$  ...  $z_{N_Z}$

$X$

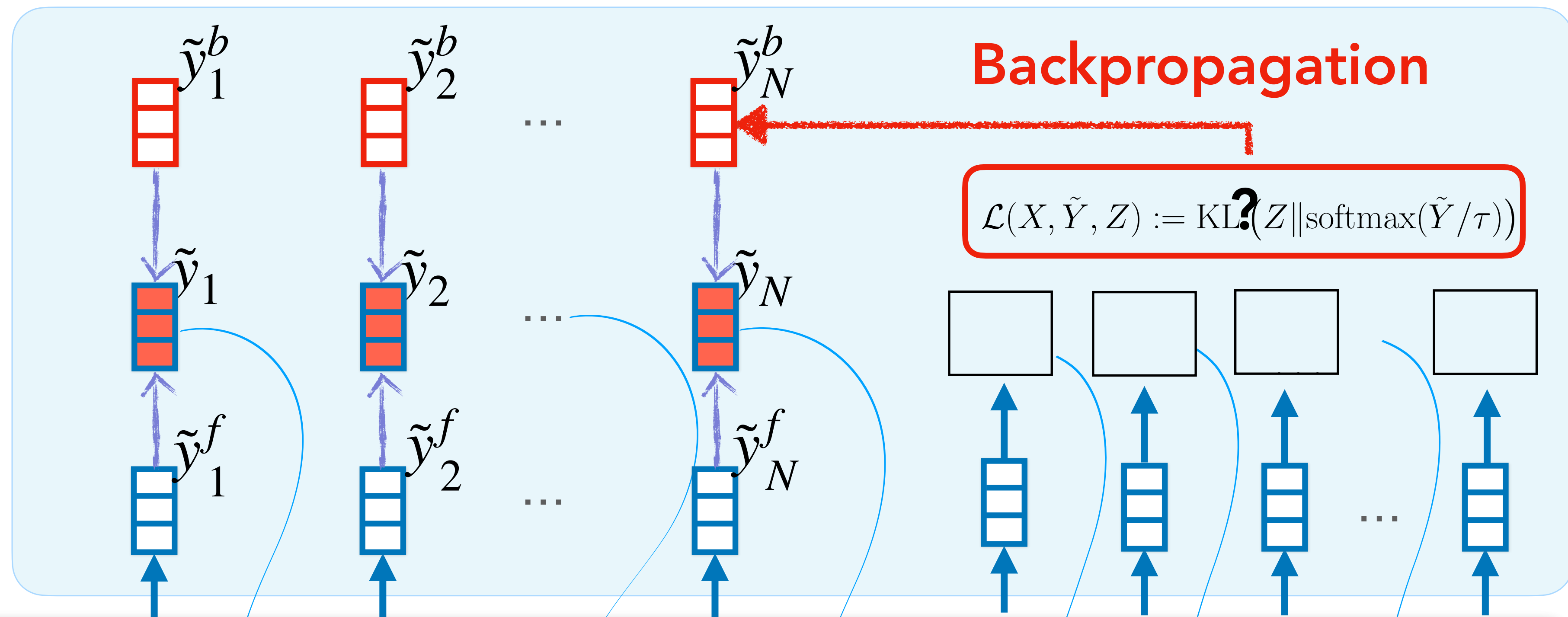
$Z$



# Counterfactual Reasoning?

Distance-Loss:  $(Y, Z)$

$Y$



LM

$x_1$   $x_2$  ...  $x_{N_X}$

$z_1$   $z_2$   $z_3$  ...  $z_{N_Z}$

$X$

$Z$

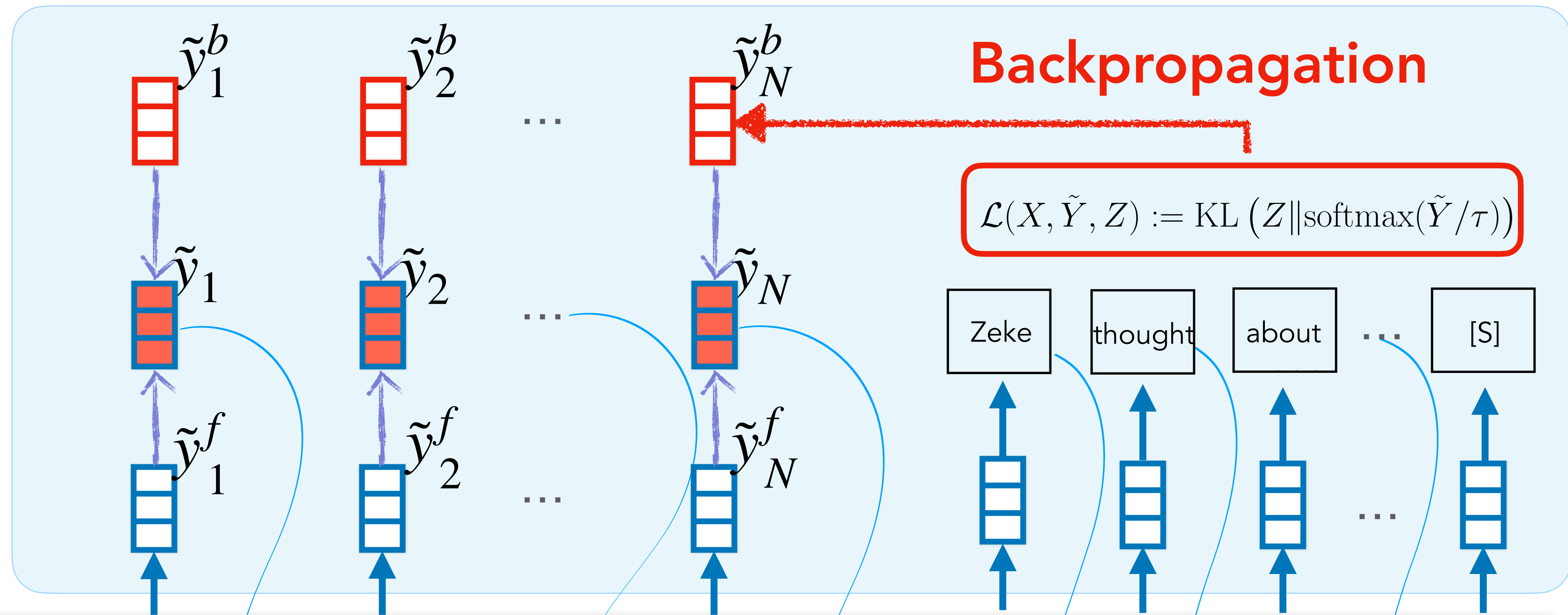




# Counterfactual Reasoning?

Distance-Loss:  $(Y, Z)$

$Y$



$x_1$   $x_2$  ...  $x_{N_X}$

Zeke was throwing a party.  
[Counterfactual] All his friends  
were dressing up for this Game  
of Thrones themed party.

$X$

$z_1$   $z_2$   $z_3$  ...  $z_{N_Z}$

Zeke thought about being  
a vampire or a wizard. Then  
he decided on a scarier  
costume. Zeke dressed up  
like a skeleton.

$Z$



# Counterfactual Reasoning?

Distance-Loss:  $(Y, Z)$

$Y$

Zeke thought about **Lannister**, but he didn't want to look like a Lannister. He wanted to look like a Stark. Zeke dressed up like a **Stark**.

Backpropagation

$$\mathcal{L}(X, \tilde{Y}, Z) := \text{KL} (Z \| \text{softmax}(\tilde{Y} / \tau))$$

LM

$x_1$   $x_2$  ...  $x_{N_X}$

Zeke was throwing a party.  
[Counterfactual] All his friends were dressing up for this Game of Thrones themed party.

$X$

$z_1$   $z_2$   $z_3$  ...  $z_{N_Z}$

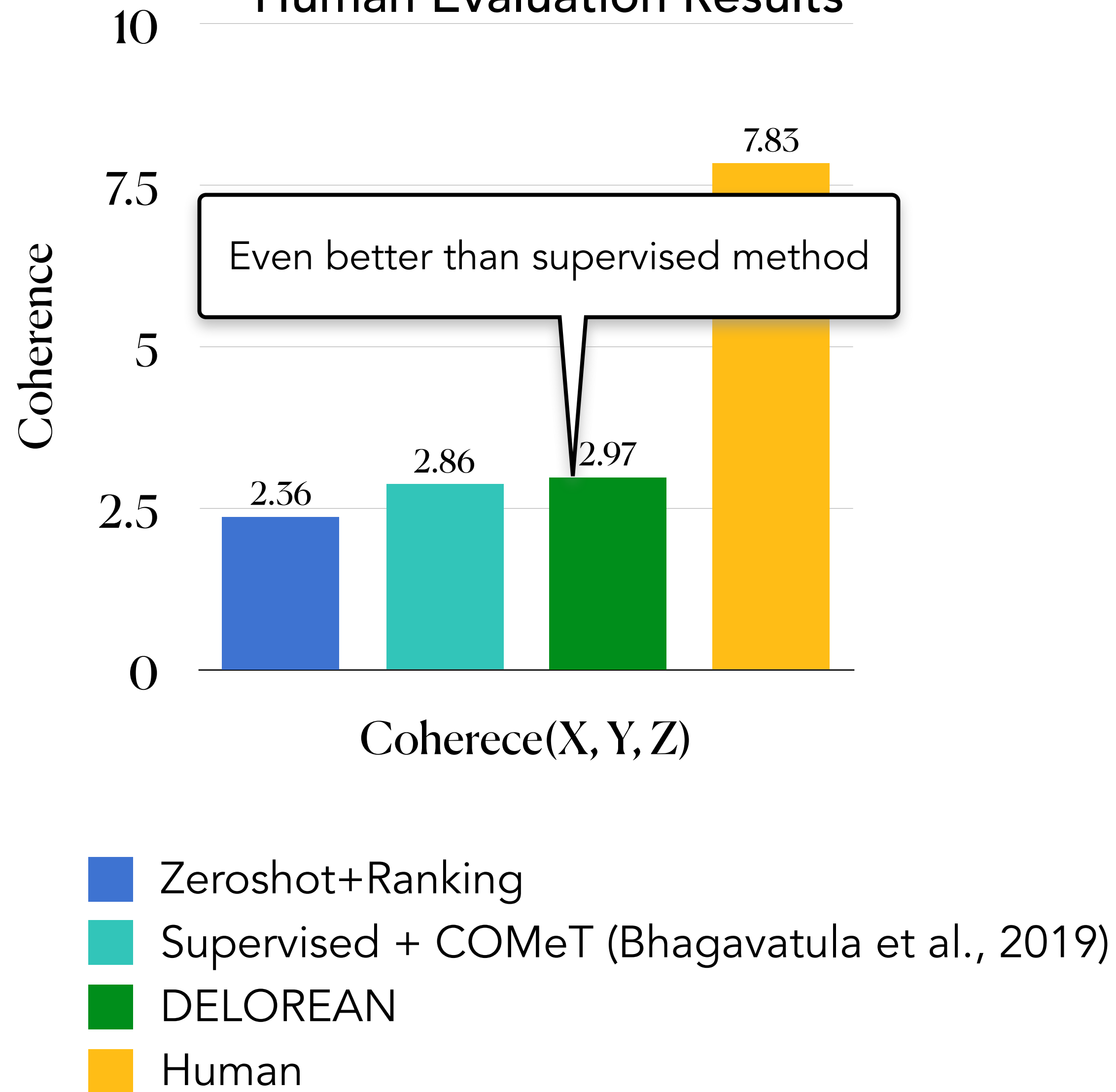
Zeke thought about being a vampire or a wizard. Then he decided on a scarier costume. Zeke dressed up like a skeleton.

$Z$

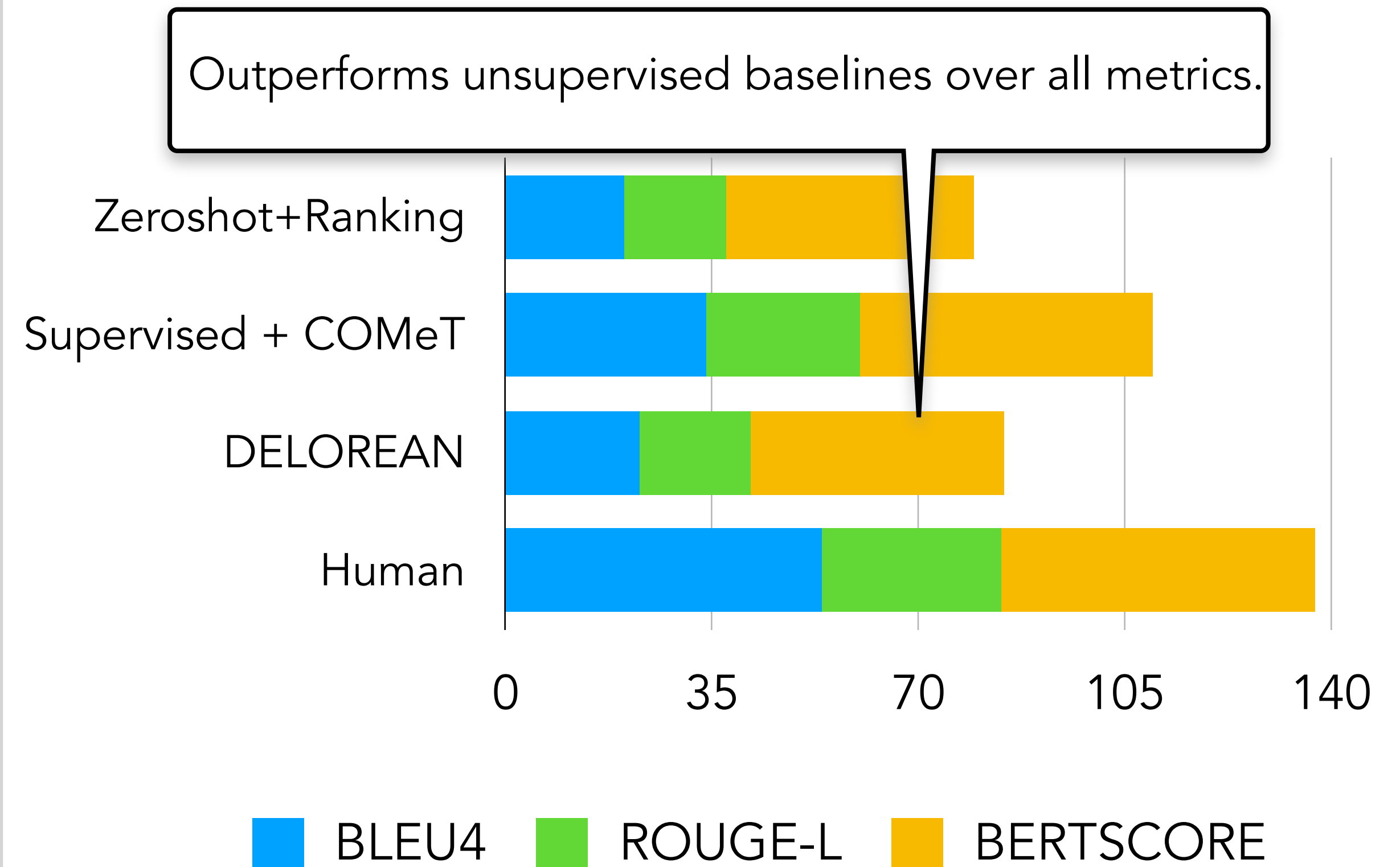


# Abductive Reasoning

## Human Evaluation Results



## Automatic Evaluation Results

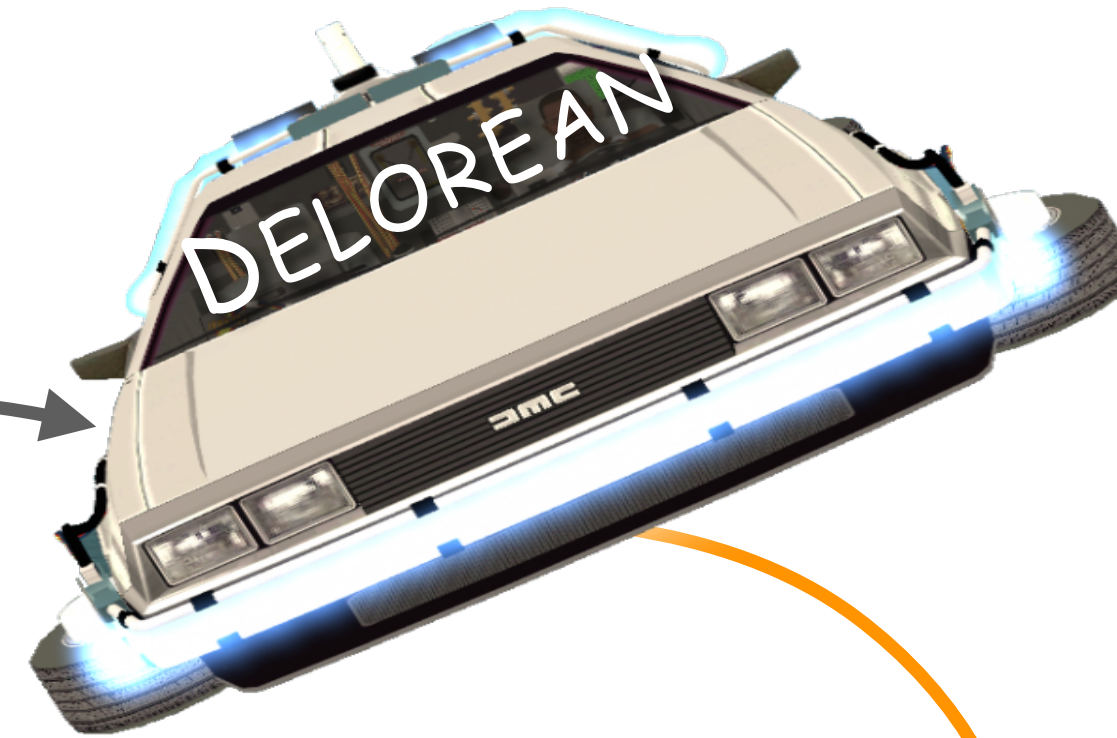


Please check the paper for more baselines ...

# Abductive Reasoning

## *Past Observation*

Ray drove his car on a steep mountain road.



## *Future Observation*

Ray was fine but his car was totaled.

## *Hypothesis*

*As he drove the car to the top of the mountain, his car is hit by a car.*



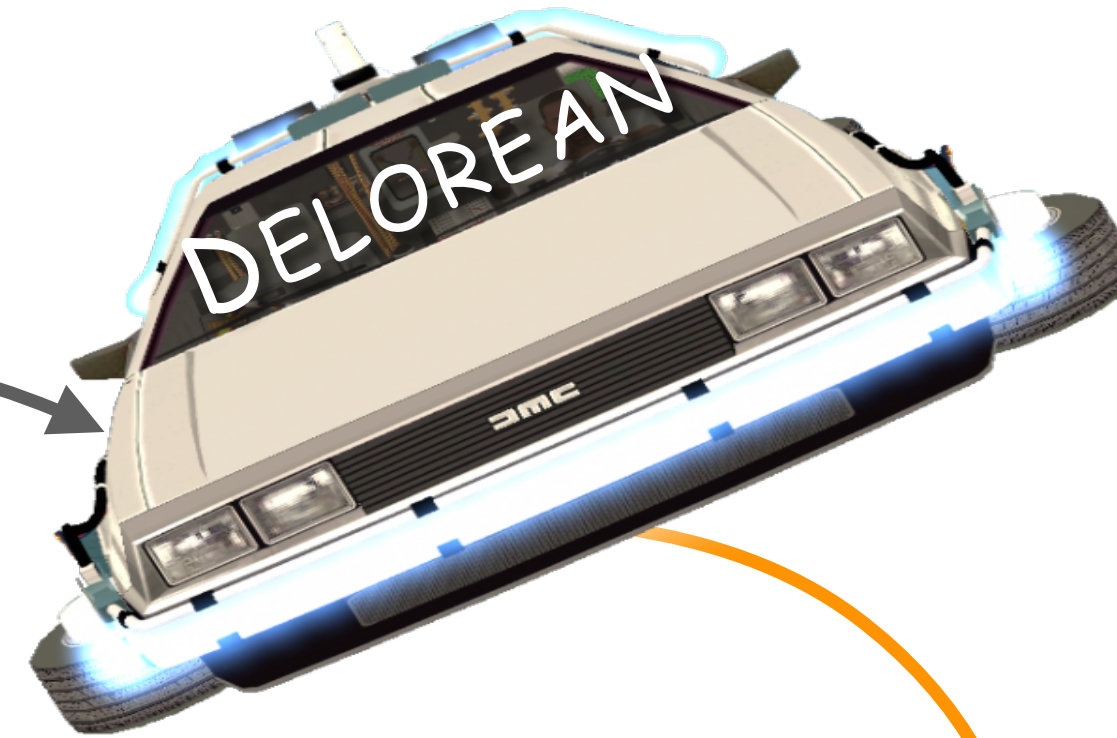
# Counterfactual Reasoning

## Story Context

Tara wanted to buy a new shirt for her upcoming school formal.

**[Original]** She went to the mall with her mom.

**[Counterfactual]** She knew of a cool place online that did custom fits really cheaply, and ordered from there.



## Original Ending

They browsed shirts from a variety of stores.

Tara picked out a floral patterned shirt that she liked best.

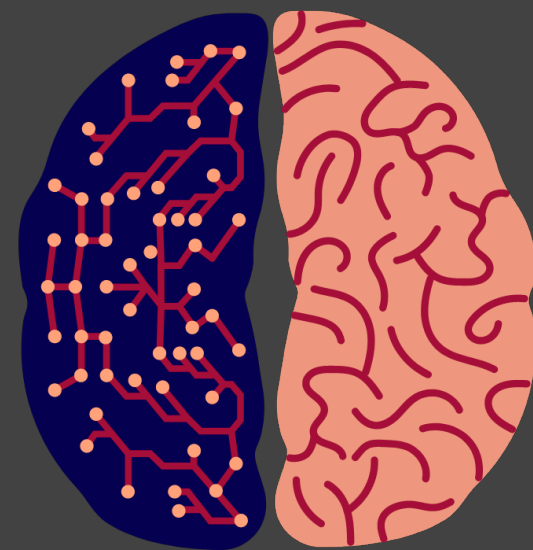
Tara looked forward to wearing it.

## Rewritten Ending

*They sent her a shirt that fit her perfectly.*

*Tara was so excited to wear it.*

*She looked forward to wearing it.*



# NEUROLOGIC DECODING

## Controlling Neural Language Generation with Logic Constraints

Ximing Lu



Peter  
West



Ronan  
LeBras



Rowan  
Zellers



Chandra  
Bhagavatula



Yejin  
Choi





# Seq2Seq

## Machine Translation

*X*

The physician told the baker that she had cancer.

*Y*

Der Arzt sagte dem Bäckerin, dass er Krebs habe.

## Dialogue Response

*X*

type	hotel
count	182
dogs allowed	don't care

*Y*

There are 182 hotels if you do not care whether dogs are allowed .

*Y*



## Language Model

## COMMONGEN

*X*

{ food, table, sit, front }

*Y*

The man sat with his food at the front of the table.

*X*



## Image Captioning

*X*



*Y*

Man in blue wetsuit is surfing on wave.

# COMMONGEN (Lin et al. EMNLP 2020)

missing keyword { **lose**, **ride** }



A man is trying to keep his **balance** as he **falls** off a **board** .

should use **all** given  
keywords ....

**Fine-tuned Language Model**

$X$

{ **board**, **lose**, **ride**, **fall**, **balance** }





# Machine Translation (Stanovsky et al., 2019)

should be female inflection **Bäckerin**

Der Arzt sagte dem Bäcker, dass er Krebs habe.

should use female  
inflection for  
women baker ....

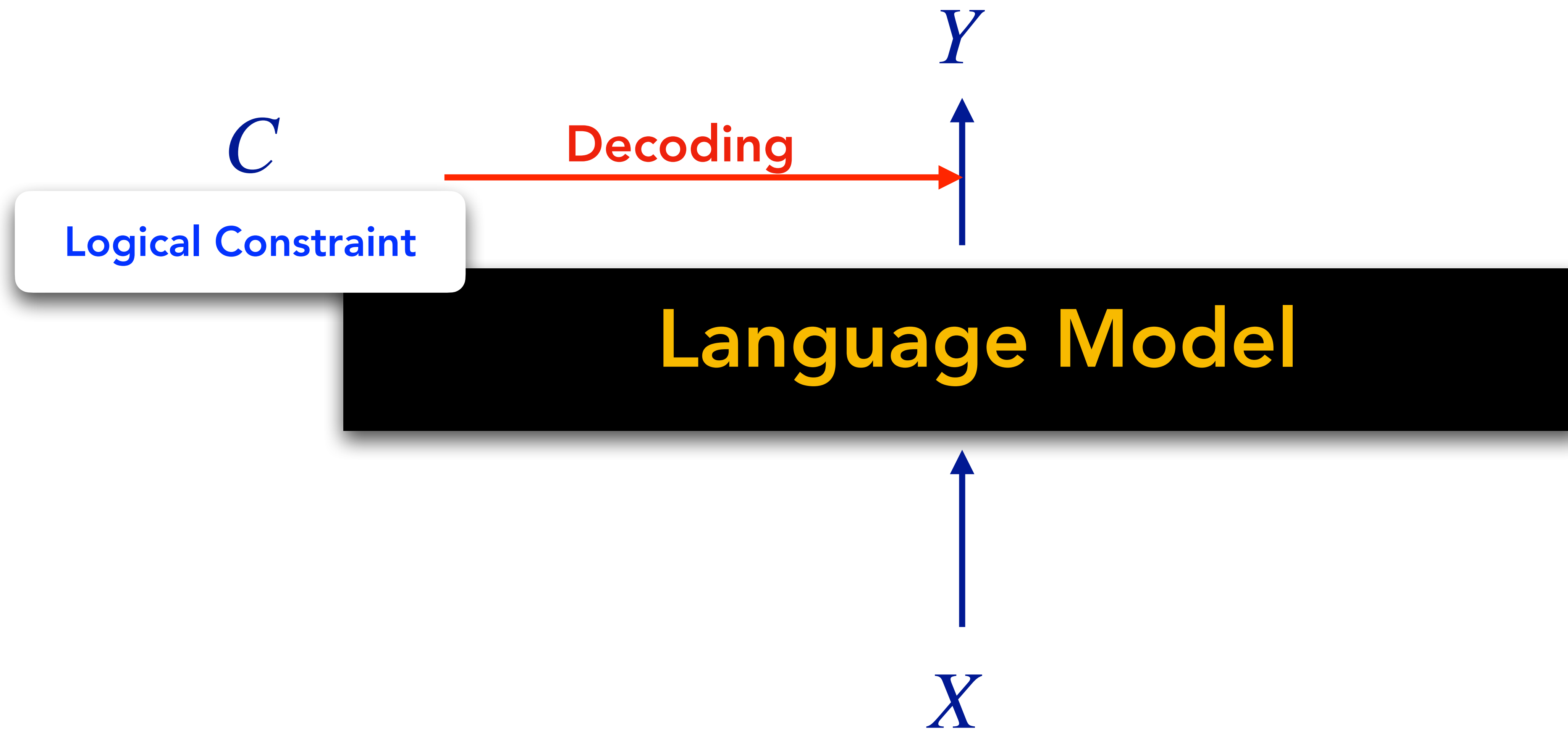
**Translation Model**

X

The physician told the baker that she had cancer.



# NEUROLOGIC DECODING





# Logical Constraints

Generate response to user's query that must contain retrieved information

$$\mathcal{D}(\text{Hilton Hotel}) \wedge \mathcal{D}(\text{address}) \wedge \mathcal{D}(\text{phone number})$$

Generate a oil-free chow mien recipe that uses  
at least one type of proteins and one type of vegetables

$$\mathcal{D}(\text{noodle}) \wedge \neg \mathcal{D}(\text{oil}) \wedge (\mathcal{D}(\text{beef}) \vee \mathcal{D}(\text{chicken}) \vee \dots) \wedge (\mathcal{D}(\text{broccoli}) \vee \mathcal{D}(\text{lettuce}) \vee \dots)$$

# Logical Constraints

CNF Form

a product of sums or an AND of ORs

$$\underbrace{(\mathcal{D}_1 \vee \mathcal{D}_2 \cdots \vee \mathcal{D}_i)}_{\mathcal{C}_1} \wedge \cdots \wedge \underbrace{(\mathcal{D}_k \vee \mathcal{D}_{k+1} \cdots \vee \mathcal{D}_l)}_{\mathcal{C}_m}$$

$\mathcal{D}_i$  : literal

$\mathcal{C}_i$  : clauses



$\forall \mathcal{C}_i, \mathcal{C}_i \text{ is true}$



$$\sum_i^m \mathcal{C}_i = m$$



# Objective

$$\hat{\mathbf{y}} = \arg \max_{\mathbf{y} \in \mathcal{Y}} P_{\theta}(\mathbf{y}|\mathbf{x})$$

subject to

some logical constraints

**CNF** 


$$\hat{\mathbf{y}} = \arg \max_{\mathbf{y} \in \mathcal{Y}} P_{\theta}(\mathbf{y}|\mathbf{x})$$

subject to

$$\sum_{i=1}^m \mathcal{C}_i = m$$

**penalty method** (Fiacco, 1976) 

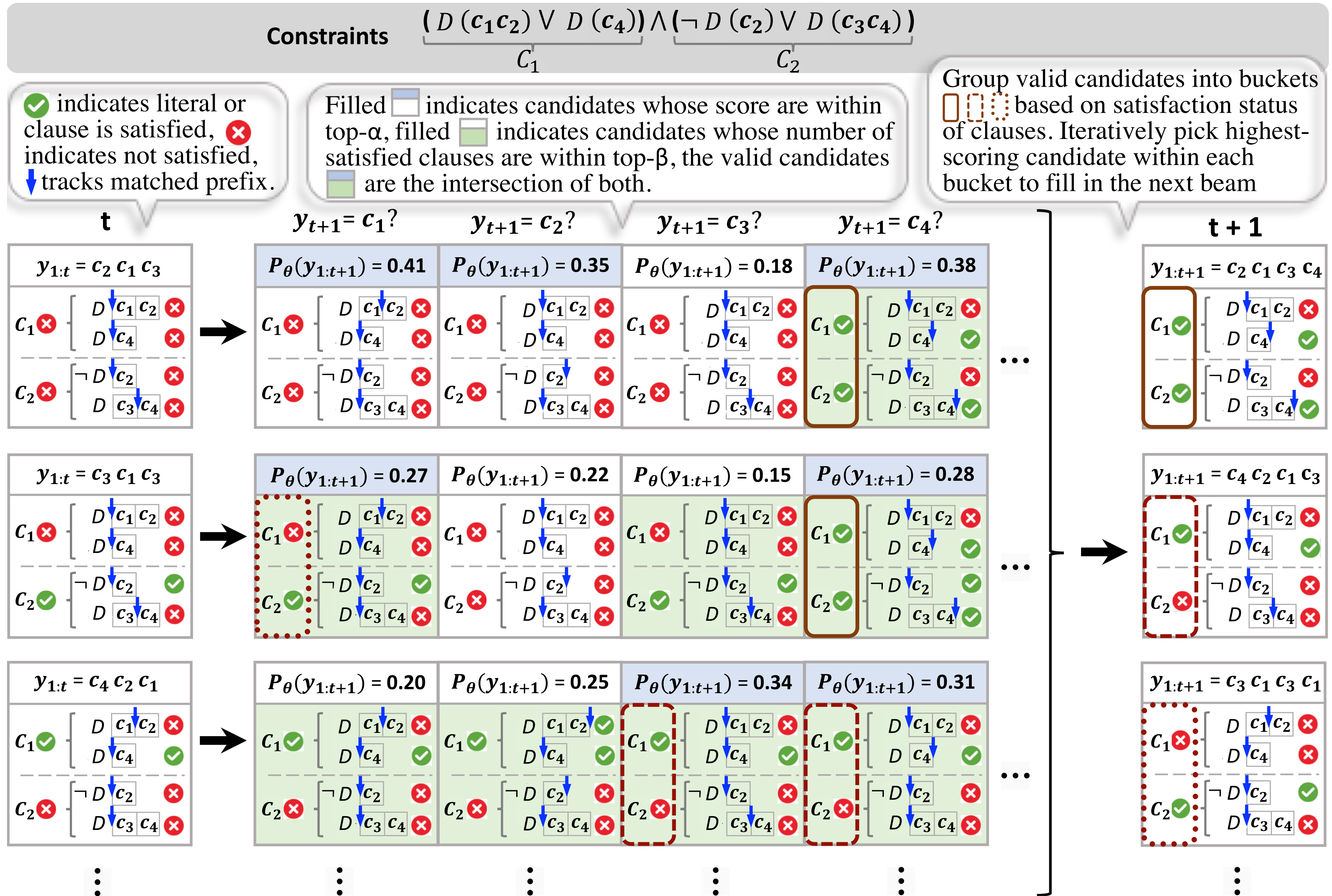
$$\hat{\mathbf{y}} = \arg \max_{\mathbf{y} \in \mathcal{Y}} P_{\theta}(\mathbf{y}|\mathbf{x}) - \alpha \left( \sum_{i=1}^m \mathcal{C}_i - m \right)^2$$

 penalty term

$$\hat{\mathbf{y}} = \arg \max_{\mathbf{y} \in \mathcal{Y}} P_{\theta}(\mathbf{y}|\mathbf{x}) + \alpha' \sum_{i=1}^m \mathcal{C}_i$$

**monotonicity** 

# Perturbed Beam Search



# COMMONGEN

$X$

{ board, lose, ride, fall, balance }

Finetuned Language Model (i.e., Supervised)

beam search

Someone **loses** his **balance**  
and **falls** off his bike .

missing keyword { **board**, **ride** }

nucleus sampling

A man is trying to keep his  
**balance** as he **falls** off a **board** .

missing keyword { **lose**, **ride** }

Nuero**Logic** decoding

Someone **loses** **balance** on the  
**ride** and **falls** off the **board** .

$C$

(**board**  $\vee$  **boards**)  $\wedge$  (**lose**  $\vee$  **loses**  
 $\vee$  **lost**)  $\wedge$  (**ride**  $\vee$  **rides**  $\vee$  **rode**  
 $\vee$  **ridden**)  $\wedge$  (**fall**  $\vee$  **falls**  $\vee$  **fell**  
 $\vee$  **fallen**)  $\wedge$  (**balance**  $\vee$  **balances**)



# COMMONGEN

$X$

{ board, lose, ride, fall, balance }

Pretrained Language Model (Unsupervised!!!)

beam search

The woman, whose name has not been released, was taken to a local hospital, where she was listed in stable condition, according to the sheriff's office.

completely irrelevant

nucleus sampling

A woman and a man have been arrested in connection with the shooting death of an unarmed black man in Ferguson, Missouri, on Aug. 12, 2014.

completely irrelevant

NueroLogic decoding

A woman **lost** her **balance** riding a horse, **falling** off the horse, and hitting her head on a **board**.

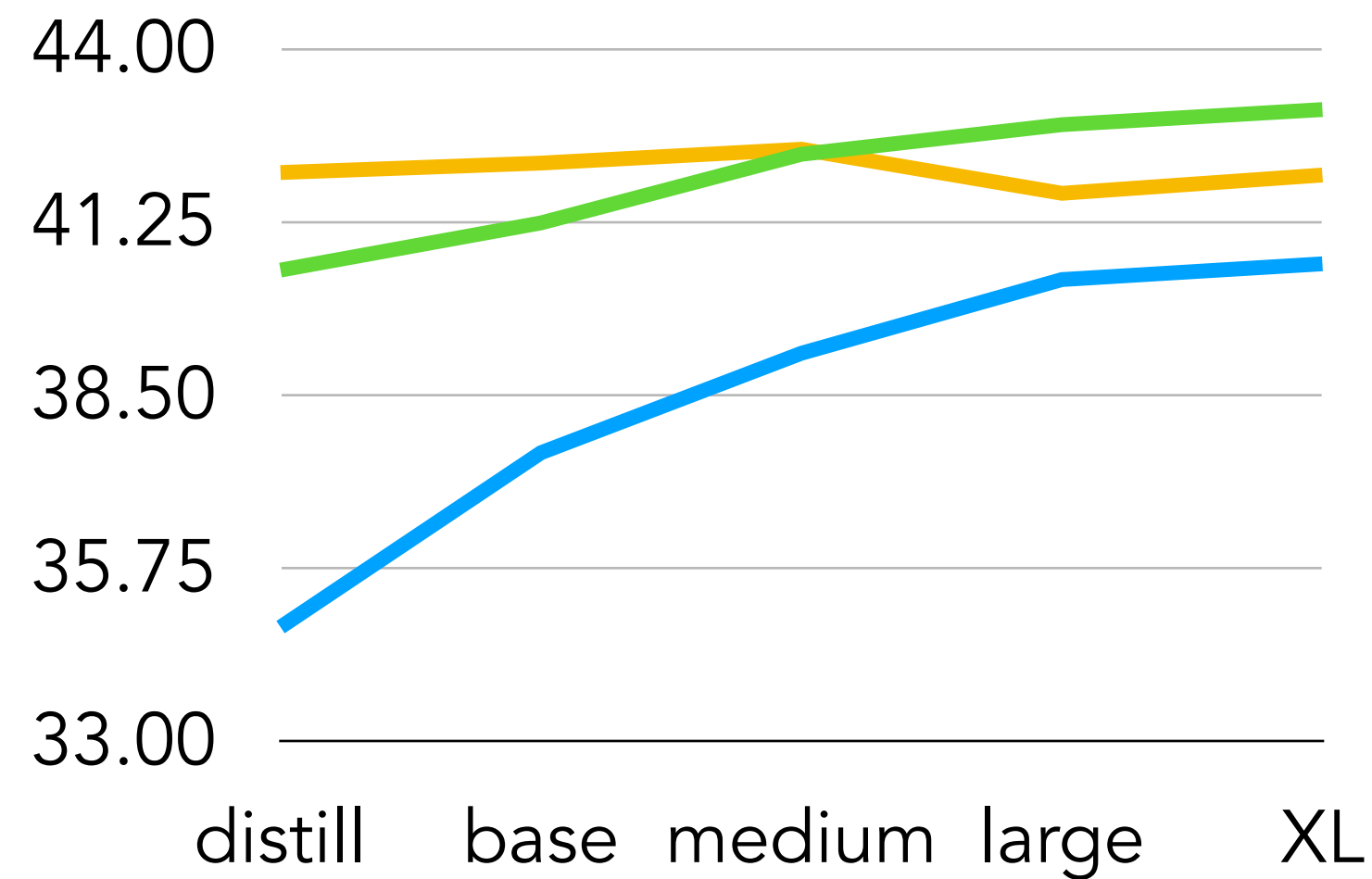
$(\text{board} \vee \text{boards}) \wedge (\text{lose} \vee \text{loses} \vee \text{lost}) \wedge (\text{ride} \vee \text{rides} \vee \text{rode} \vee \text{ridden}) \wedge (\text{fall} \vee \text{falls} \vee \text{fell} \vee \text{fallen}) \wedge (\text{balance} \vee \text{balances})$

$C$

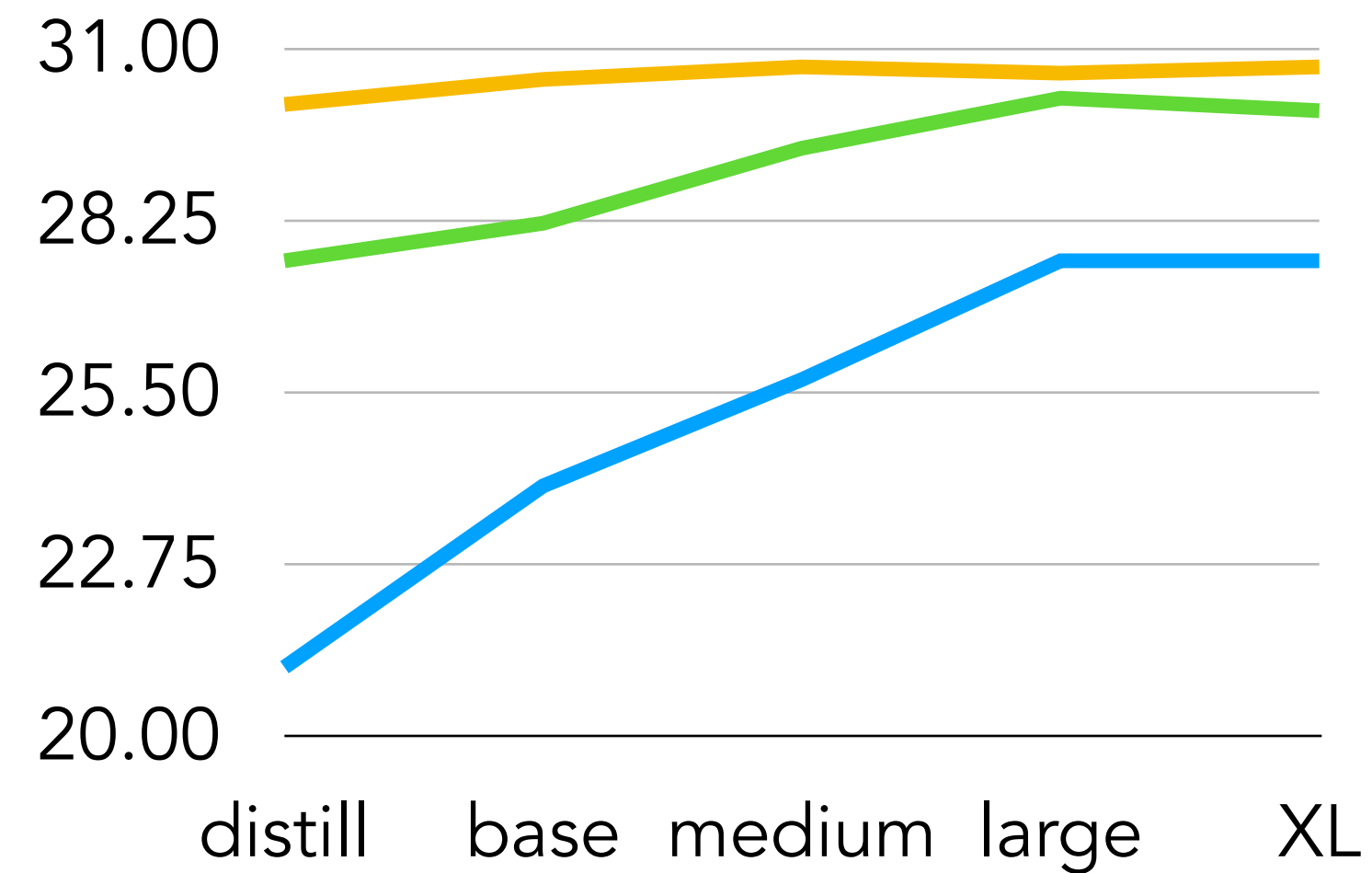
# COMMONGEN (Zero-shot)

- beam search (supervised)
- NeuroLogic (supervised)
- NeuroLogic (zero-shot)

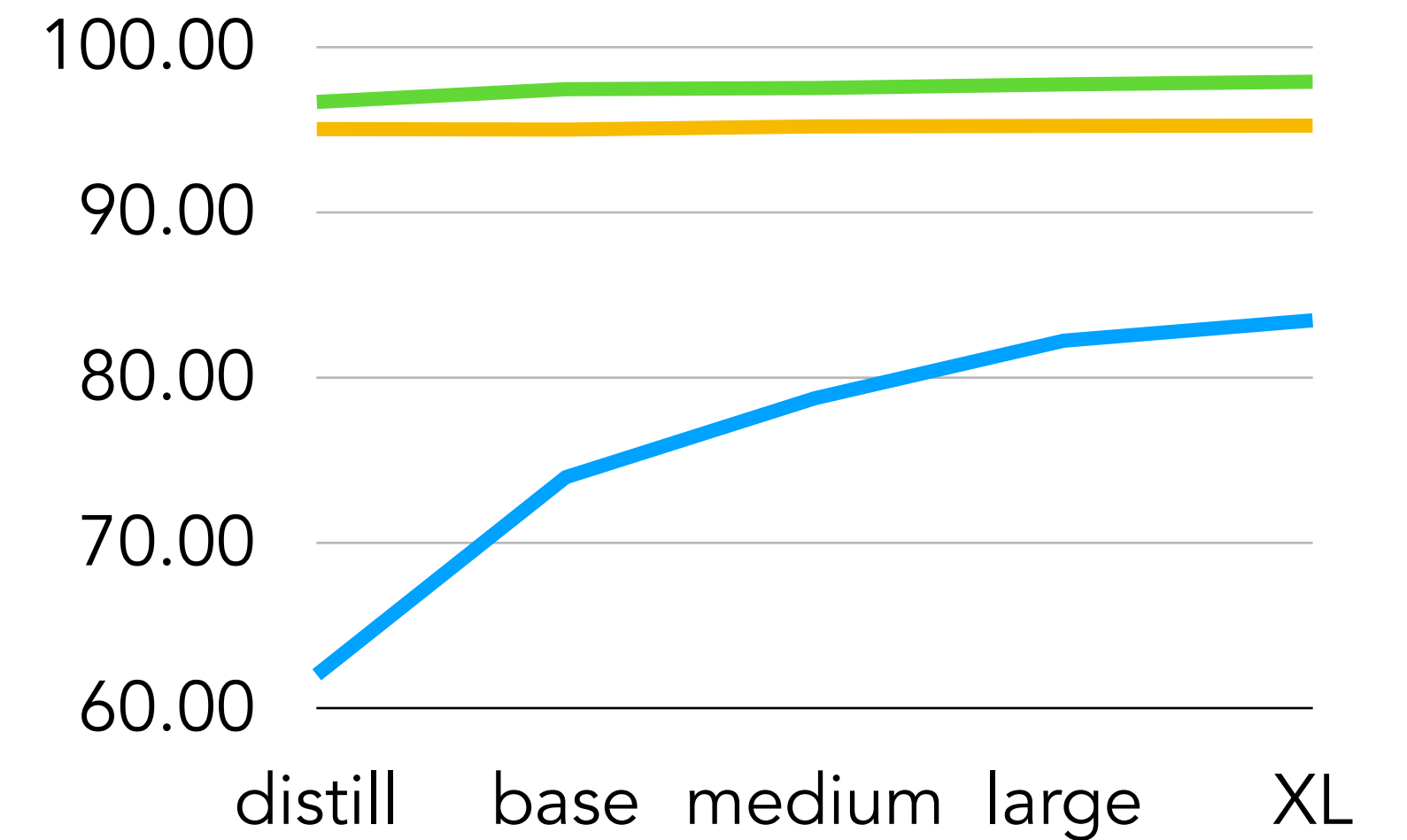
## ROUGE-L



## METEOR



## Coverage



Unsupervised NeuroLogic  
outperforms  
supervised approaches

Unsupervised NeuroLogic on **smaller**  
**networks** outperforms  
supervised approaches on **larger** networks!

# Recipe Generation

*X*

Garlic-Butter Steak

{ 2 tsp butter, 1 beef steak, 1/4 tsp soy sauce, 1 tsp parsley, 1/8 tsp salt , 1/2 tsp garlic }



Finetuned Language Model

missing ingredient { soy sauce }  
unprovided ingredient { pepper }

beam search

NueroLogic decoding

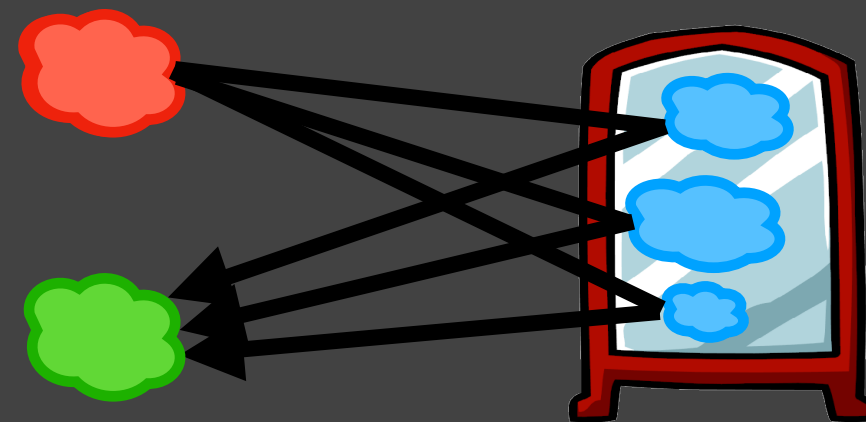
*C*

Preheat a grill for high heat. In a bowl, mix together the butter, garlic, parsley, salt and pepper. Place the steak on the hot grill and cook for 3 to 4 minute per side. Transfer the steak to a cutting board and cut it into 1 inch wide strips.

Mix 1 tablespoon butter, parsley, garlic and soy sauce. Sprinkle steak with salt. In a large skillet, heat remaining butter over medium heat. Add steak; cook until meat reaches desired doneness, 4-7 minutes per side. Serve with garlic butter.

butter  $\wedge$  (beef  $\vee$  steak  $\vee$  meat)  
 $\wedge$  soy sauce  $\wedge$  (parsley  $\vee$  herb)  
 $\wedge$  salt  $\wedge$  (garlic  $\vee$  vegetable)  
 $\wedge$  ( $\neg$ pork  $\wedge$   $\neg$ bean  $\wedge$   $\neg$ ...)





# Reflective Decoding: Unsupervised Paraphrasing and Abductive Reasoning



Peter West



Ximing Lu



Ari Holtzman



Chandra  
Bhagavatula



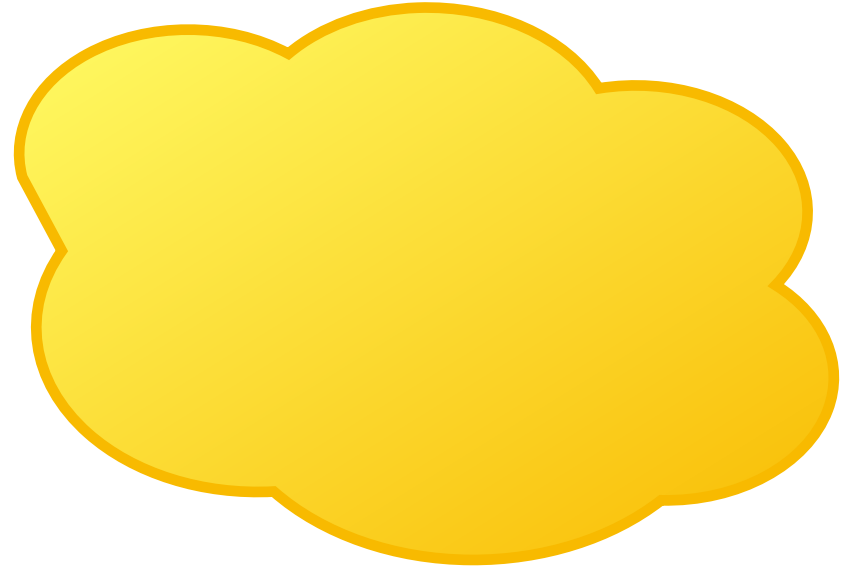
Jena Hwang



Yejin Choi

# Paraphrasing with Reflective Decoding

Input



Context





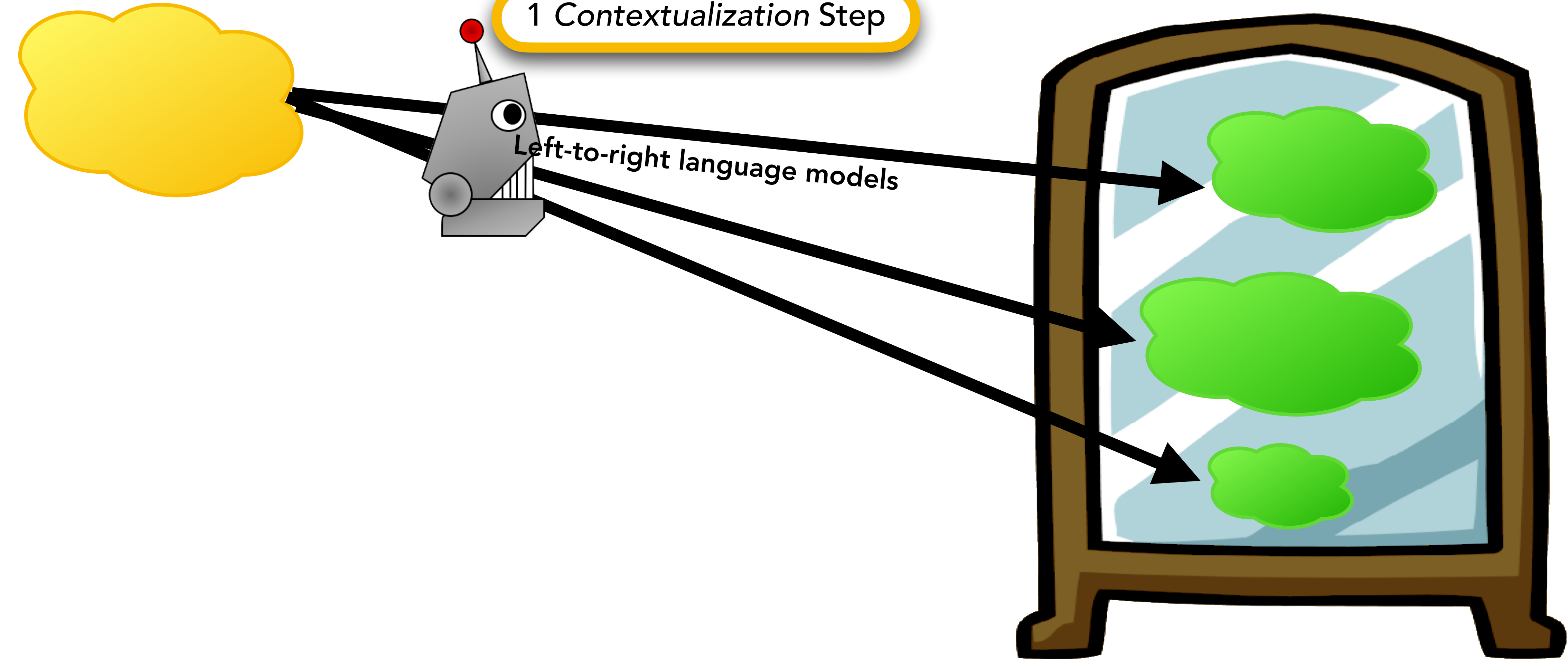
# Paraphrasing with Reflective Decoding

Input

Context

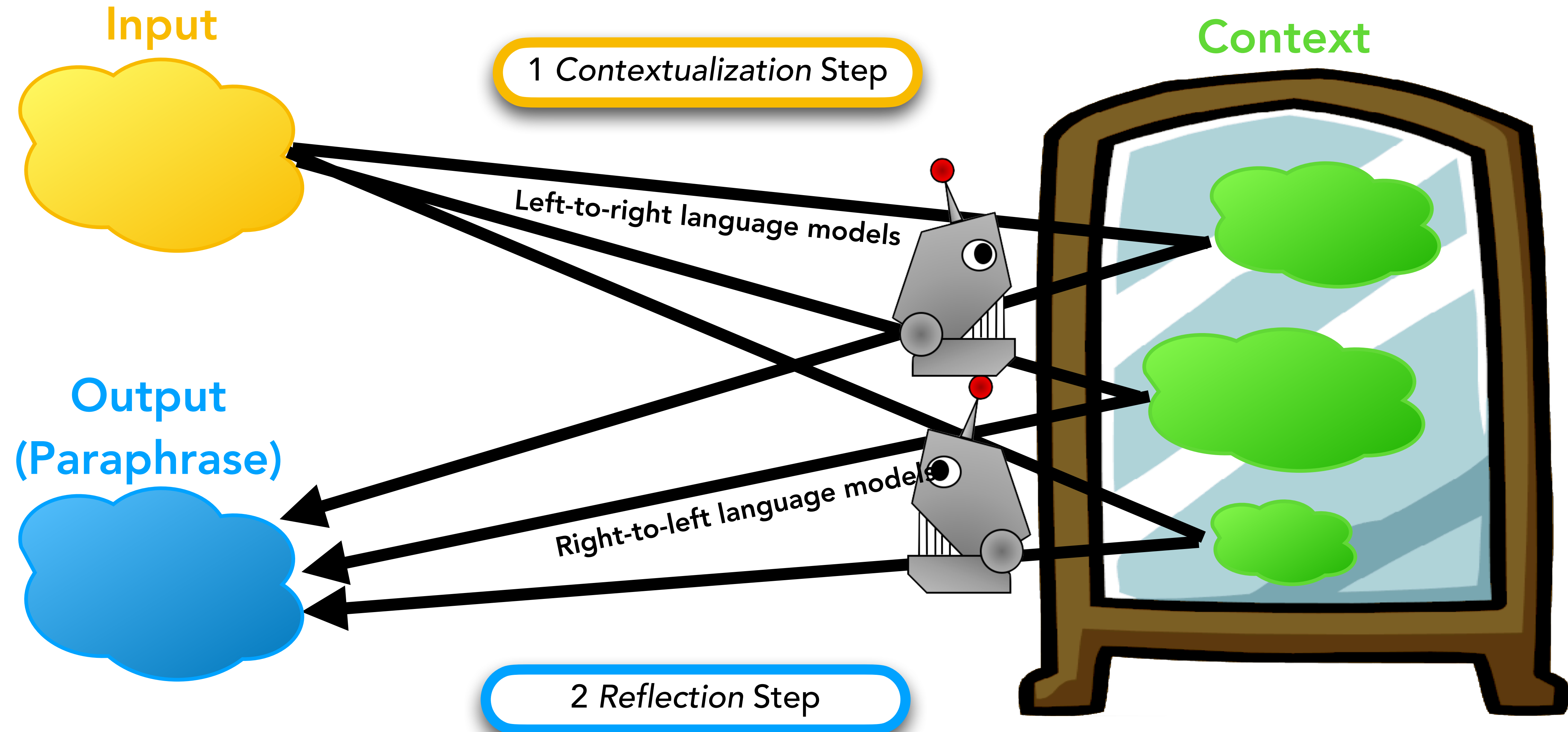
1 Contextualization Step

Left-to-right language models





# Paraphrasing with Reflective Decoding

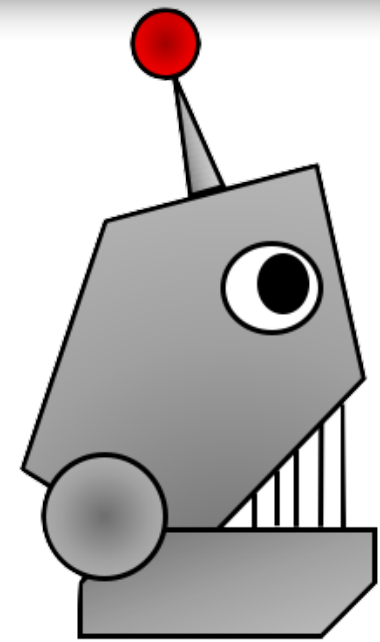


# Paraphrasing with Reflective Decoding

1 Contextualization Step

Meaning of { **Input** }

What do people  
think of  
Americans?



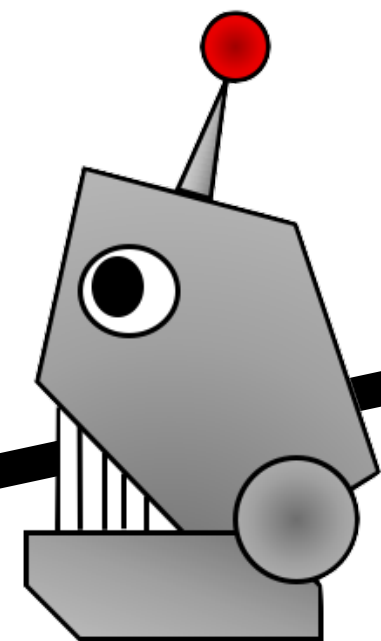
Meaning of { **Context** }

Left-to-right language models

As an American, I'm not  
qualified to answer this

Meaning of { **Output  
(Paraphrase)** }

How are  
Americans  
perceived?



Right-to-left language models

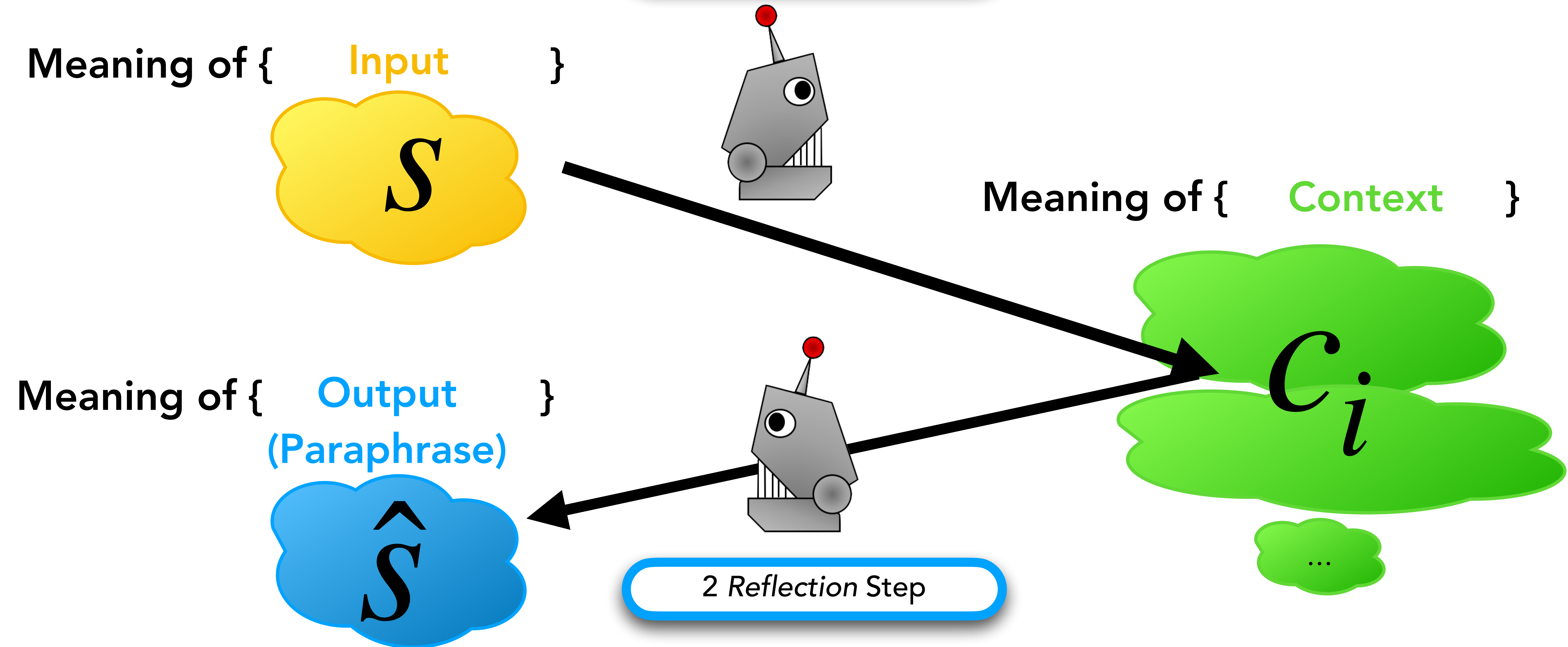
People have mixed feelings  
about Americans

...

2 Reflection Step

# Paraphrasing with Reflective Decoding

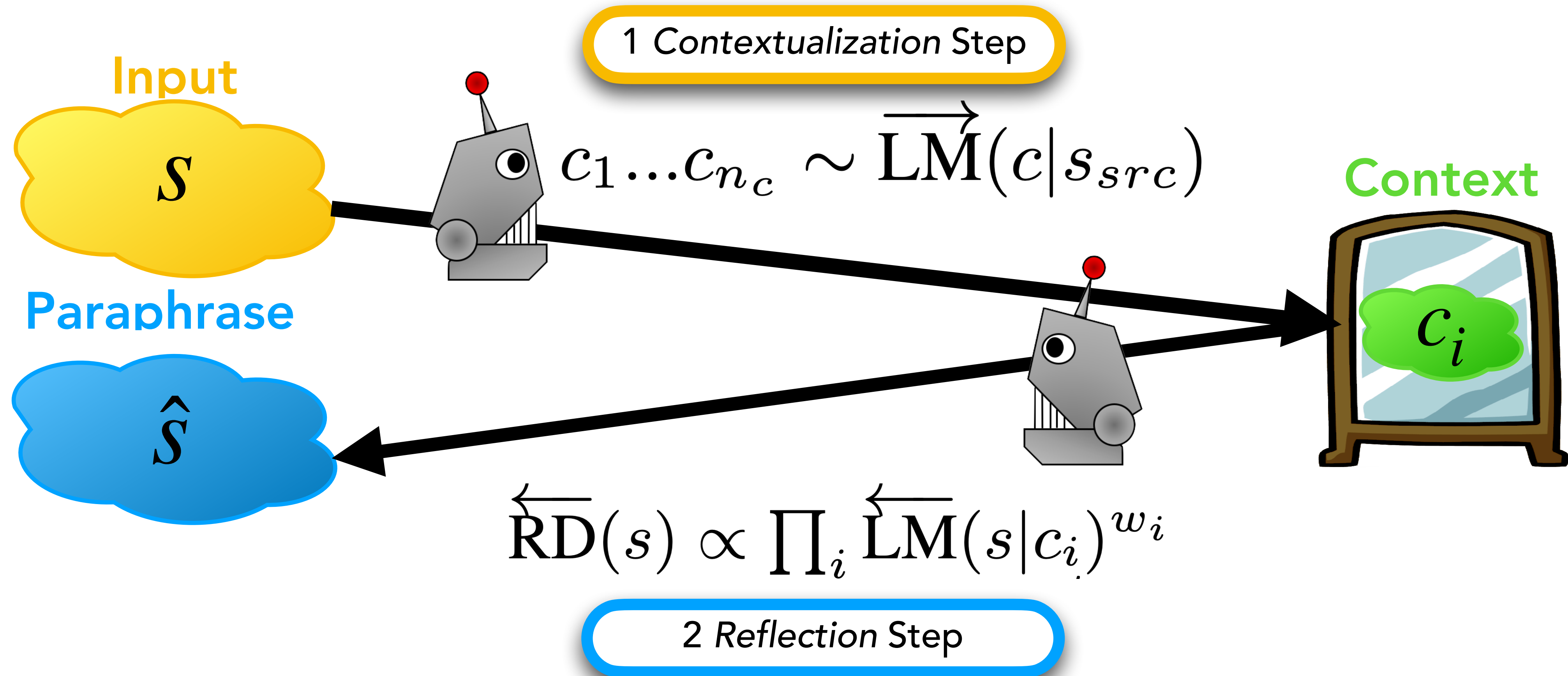
1 Contextualization Step



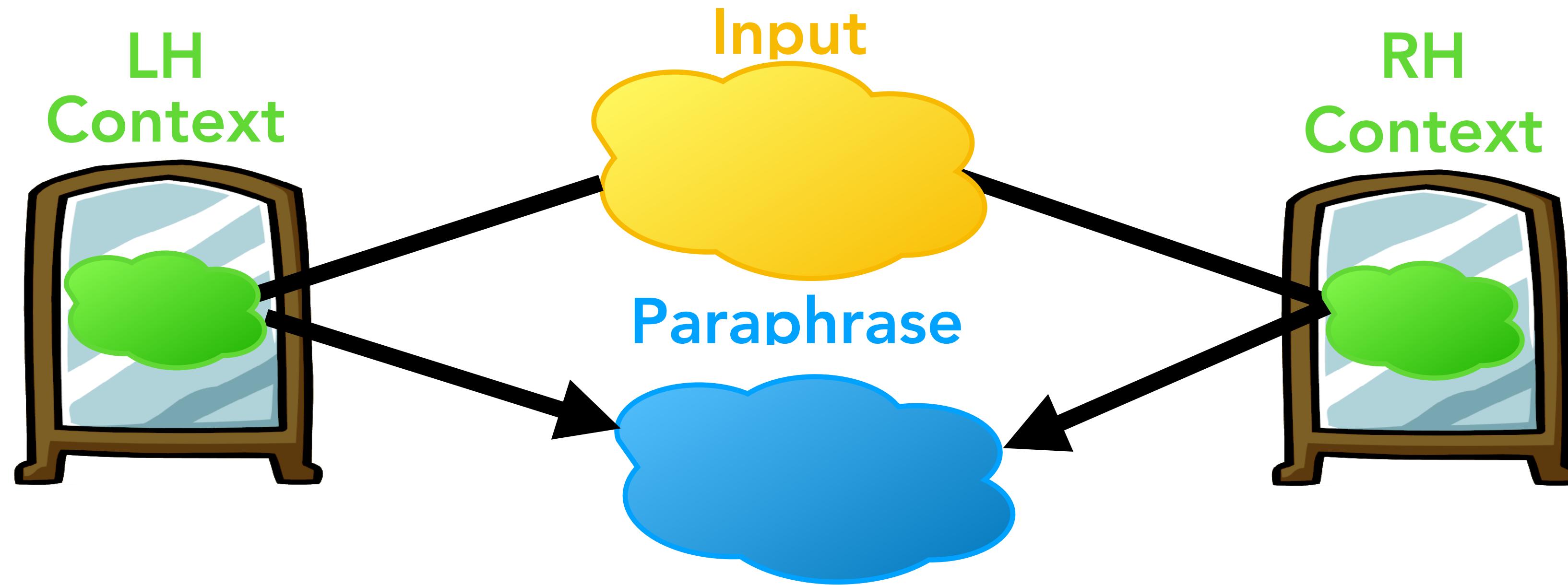


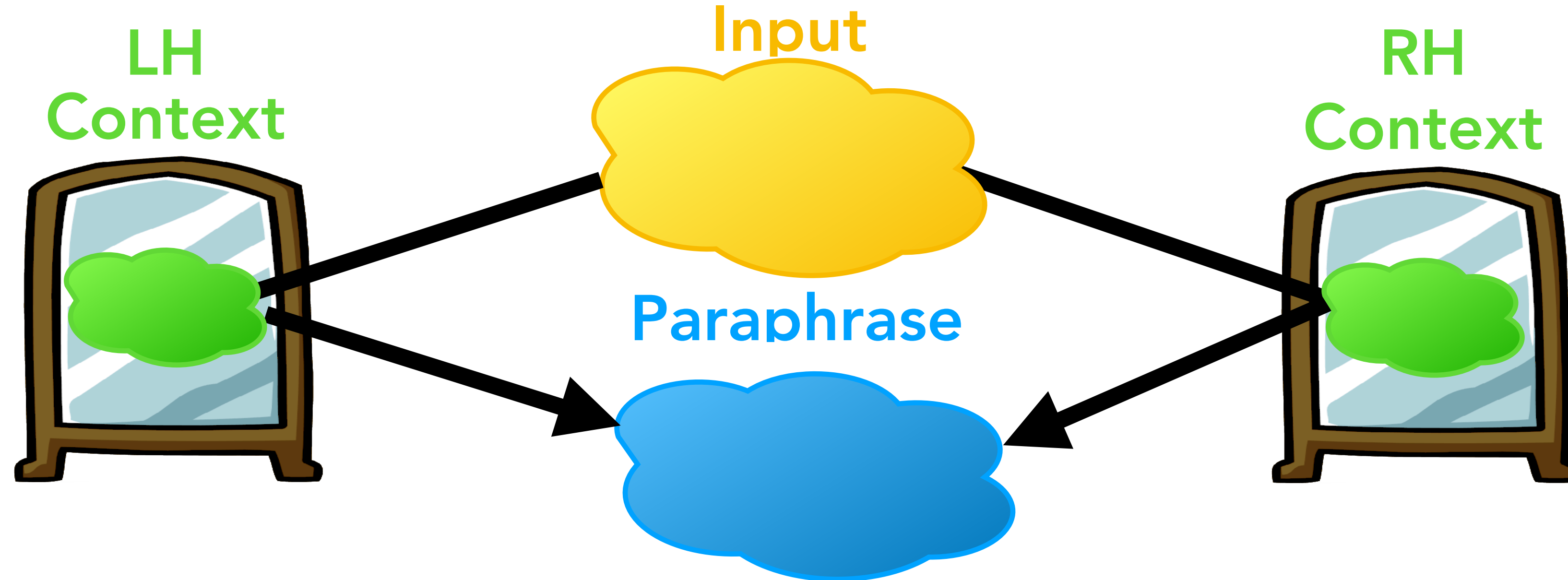
# Paraphrasing with Reflective Decoding

$$D_{KL}(\overrightarrow{\text{LM}}(c|s_{src}), \overrightarrow{\text{LM}}(c|\hat{s}))$$

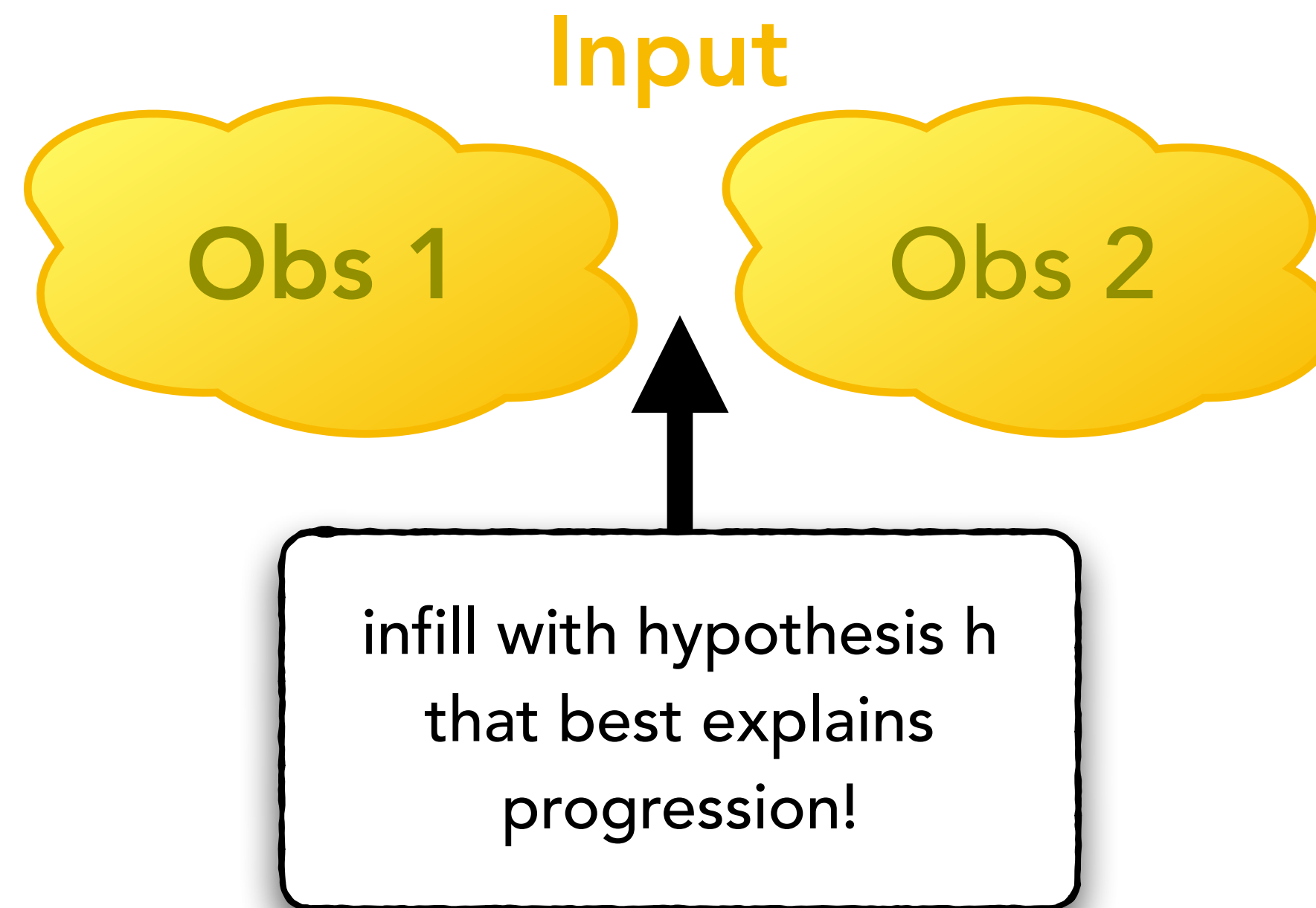


# Use case 1: paraphrasing

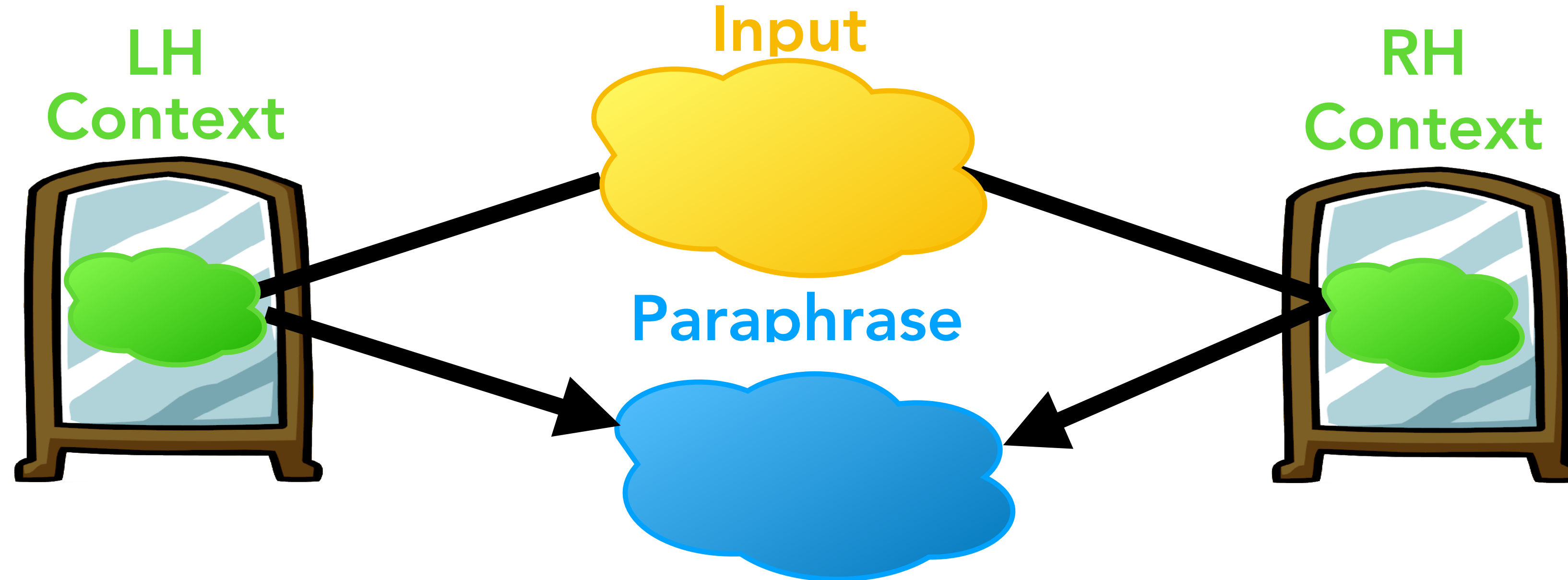




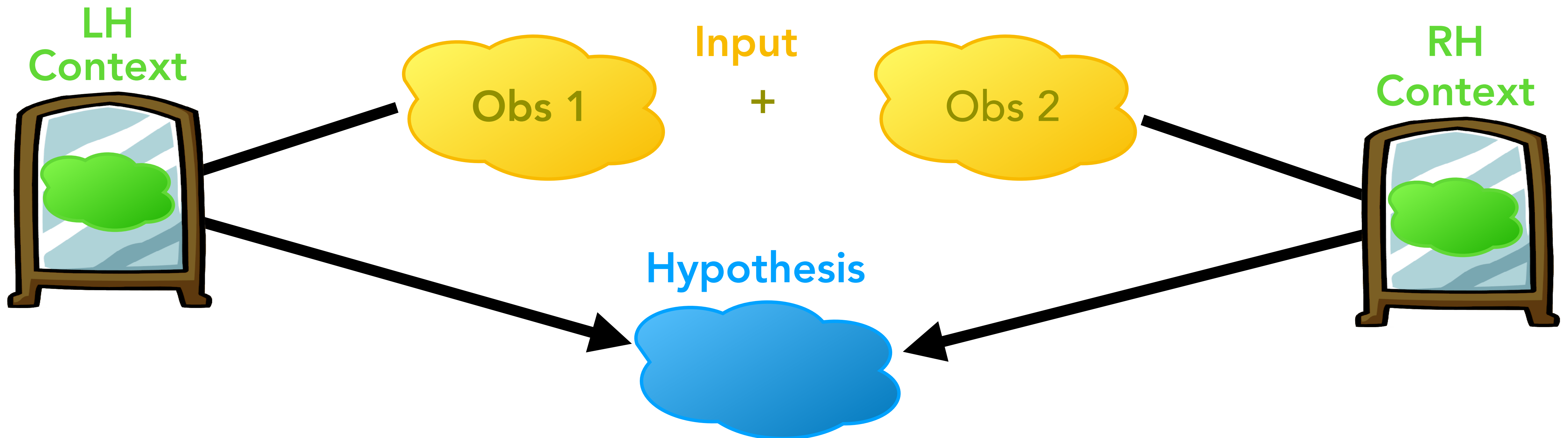
## Use case 2: Abductive NLG



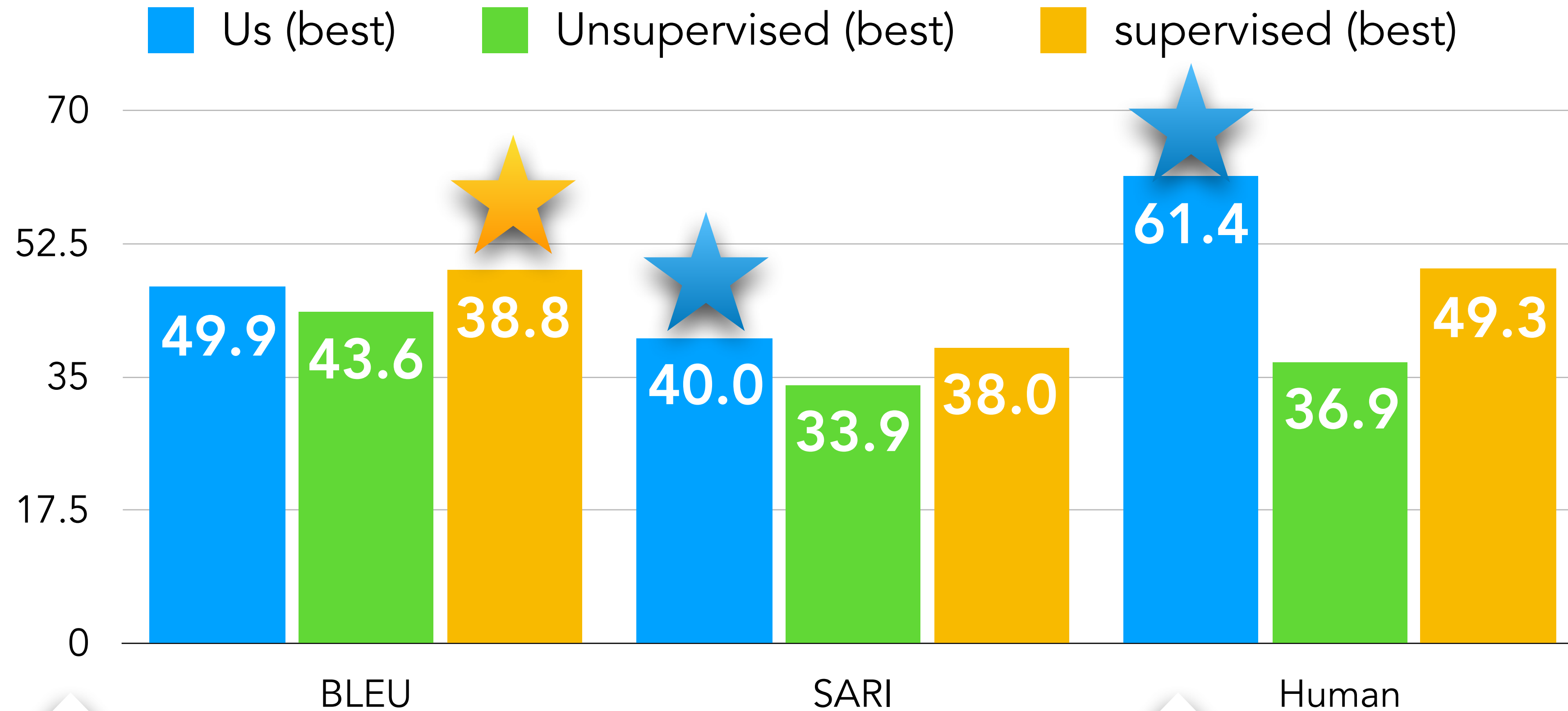




## Use case 2: Abductive NLG



# Evaluation on Paraphrasing



BLEU can be gamed by copy-and-paste the source phrase as is (without paraphrasing)  
SARI is a new major that accounts for novelty!

**Unsupervised Reflective Decoding** (not even fine-tuned on the target domain) outperforms **supervised (BART)** based on **SARI & Human Eval**

# In this talk: Reasoning as Generation

- **Part 1:** unsupervised inference-time algorithms

Reasoning thru  
**Neural Backpropagation**


DeLorean

Reasoning thru  
**Search with Logical Constraints**

NeuroLogic

Reasoning thru  
**Distributional Neural Imagination**

Reflective Decoding

=> How to make  from (off-the-shelf) neural language models



# Current Paradigm of Deep Learning

To achieve (super) human level performances on language (and vision) leaderboards

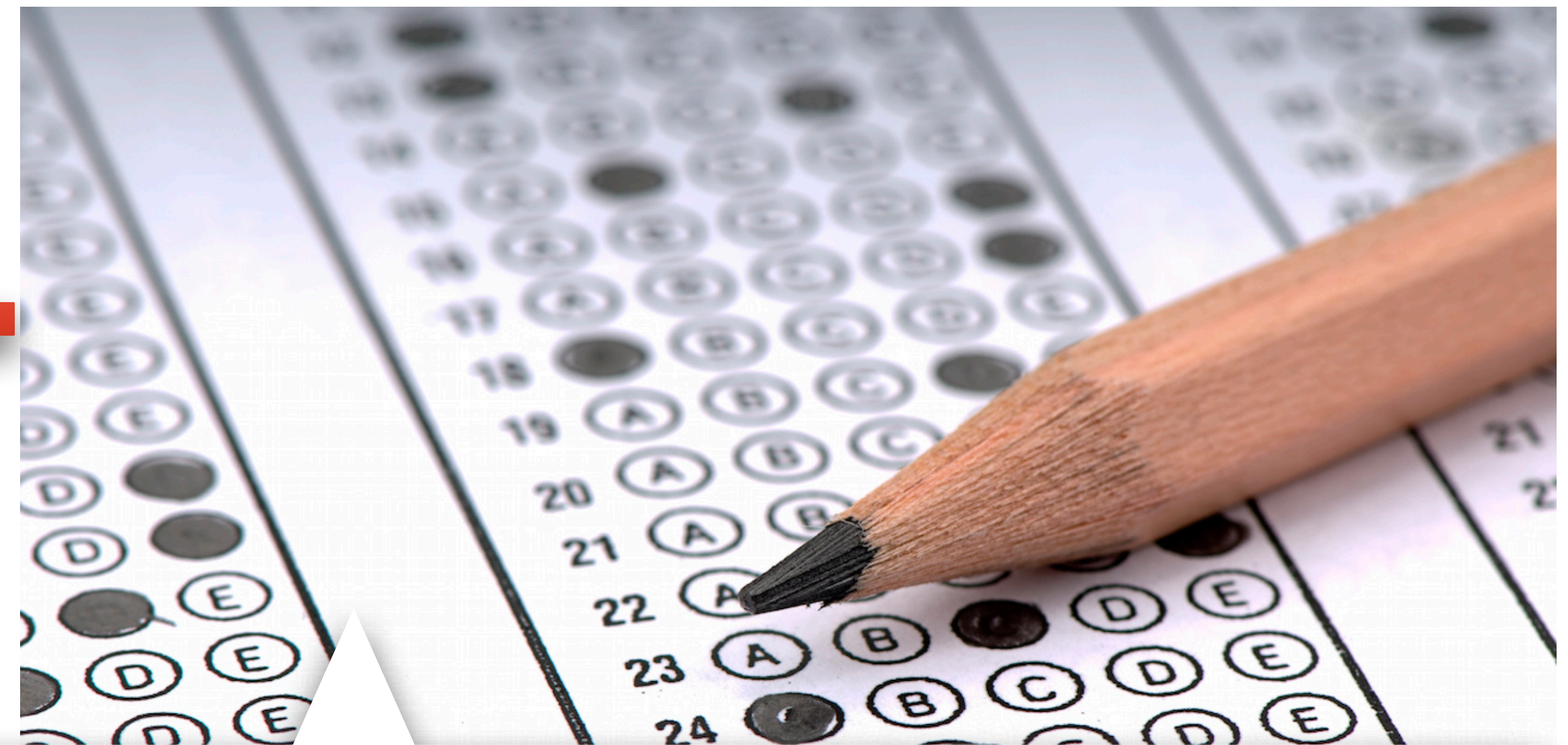
**Self-supervision** on a lot of raw text

...rict reports indicate somewhat stronger  
...early January than at the time of the  
...entered in the retail and industrial sec  
...s , that the national economy gained  
...ed , manufacturing activity [MASK] to  
...ant and equipment .

...al economic activity on balance in Dec  
...ports in November , with much of the  
...It would appear , on the basis of th



**Supervision** on a lot of exam problems



This recipe fails on generative evaluation such as  
abductive reasoning, counterfactual story revision, and  
various constrained text generation tasks with logical constraints

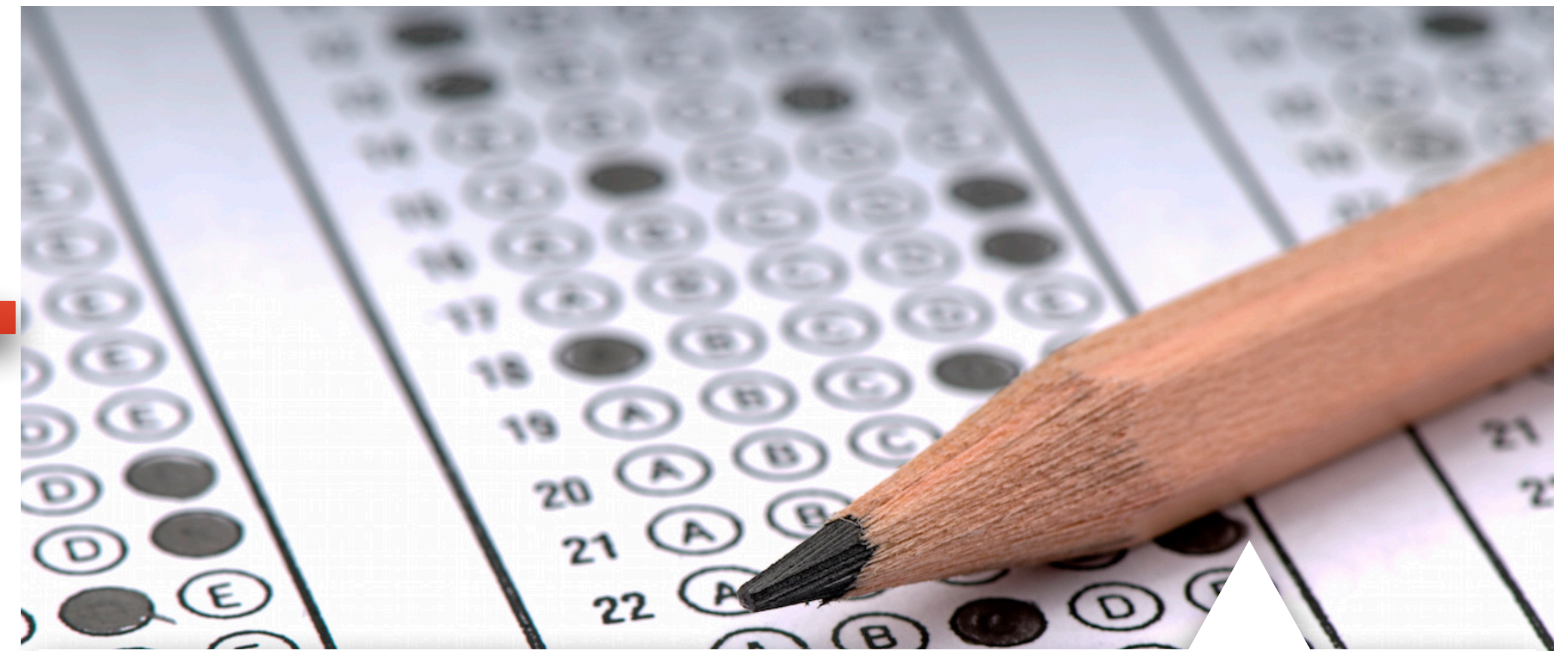
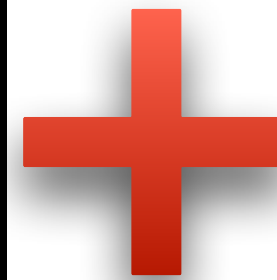


# Imagine taking a deep learning class in which...

**Self-supervision** on a lot of DL code

**Supervision** on a lot of exam problems

```
86 self.names = names
87 self.name2index = dict(zip(names, range(len(names))))
88
89
90 def __del__(self):
91     # free memory created by C to avoid memory leak
92     if hasattr(self, '__createfrom__') and self.__createfrom__ == 'C':
93         if pointer(self) is not None:
94             libbigfile.free_file(pointer(self))
95
96 def read(self, requested, isname=True):
97     if isname:
98         index_name_array = [(self.name2index[x], x) for x in requested]
99     else:
100         assert(min(requested)>=0)
101         assert(max(requested)<len(self.names))
102         index_name_array = [(x, self.names[x]) for x in requested]
103         index_name_array.sort()
104
105     npoints = len(index_name_array)
106     c_index = (c_ulonglong * npoints)()
107     for i in range(npoints):
108         c_index[i] = index_name_array[i][0]
109
110     size = self.ndims * npoints
111     pdata = (c_float * size)()
112     res = libbigfile.seq_read_memory(self, npoints, c_index, pdata)
113     assert(res)
```



NNs latch on spurious correlations & unwanted dataset biases

Can't learn concepts well enough; need to learn from declarative knowledge



# In this talk: Reasoning as Generation

- **Part 1:** unsupervised inference-time algorithms

Reasoning thru  
**Neural Backpropagation**

DeLorean

Reasoning thru  
**Search with Logical Constraints**

NeuroLogic

Reasoning thru  
**Distributional Neural Imagination**

Reflective Decoding

- **Part 2:** supervision with declarative knowledge for knowledge modeling

**Neural & Symbolic  
Commonsense Knowledge**

COMET & ATOMIC 2020

**Visually Grounded  
Commonsense Knowledge**

Visual COMET

**Social, Ethical, Moral Norms**

Social Chemistry 101

- **Part 3:** benchmarks and algorithmic bias reduction



# (COMET-) ATOMIC<sub>20</sub><sup>20</sup>:

## On Symbolic and Neural Commonsense Knowledge Graphs

— wait, doesn't GPT-3 know everything? —

To appear at AAAI 2021

Jena  
Hwang



Chandra  
Bhagavatula



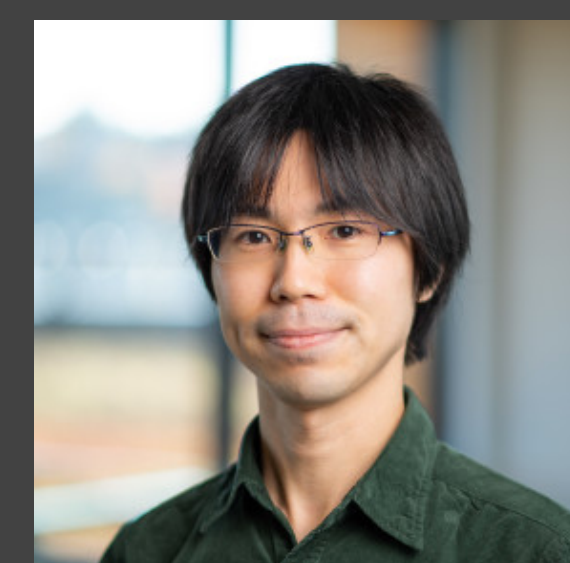
Ronan  
Le Bras



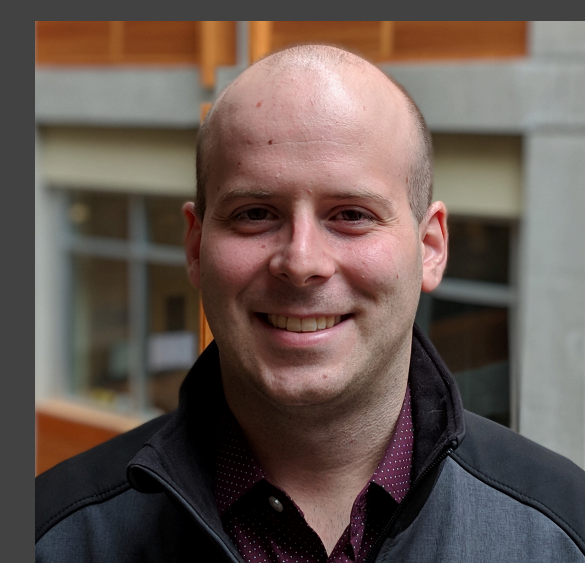
Jeff  
Da



Keisuke  
Sakaguchi



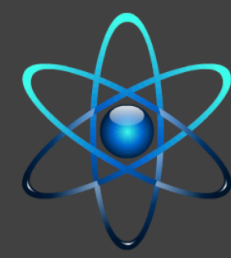
Antoine  
Bosseult



Me







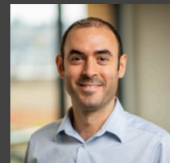
# ATOMIC: An Atlas of Machine Commonsense for If-Then Reasoning

AAAI 2019

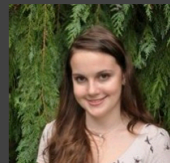
Maarten Sap



Ronan LeBras



Emily Allaway



Chandra Bhagavatula



Nicholas Lourie



Hannah Rashkin



Brendan Roof



Noah Smith



Me

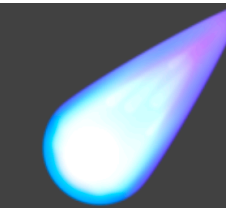


## Therapy Chabot

Kearns et al. 2020 @ CHI EA 2020

## Automated Storytelling

Ammanabrolu et al. 2020 @ arXiv:2009.00829



# COMET: Commonsense Transformers for Automatic Knowledge Graph Construction

ACL 2019

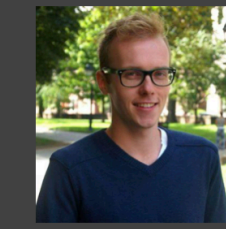
Antoine Bosselut



Hannah Rashkin



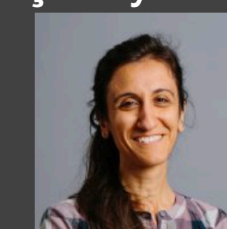
Maarten Sap



Chaitanya Malaviya



Asli Çelikyilmaz



Me



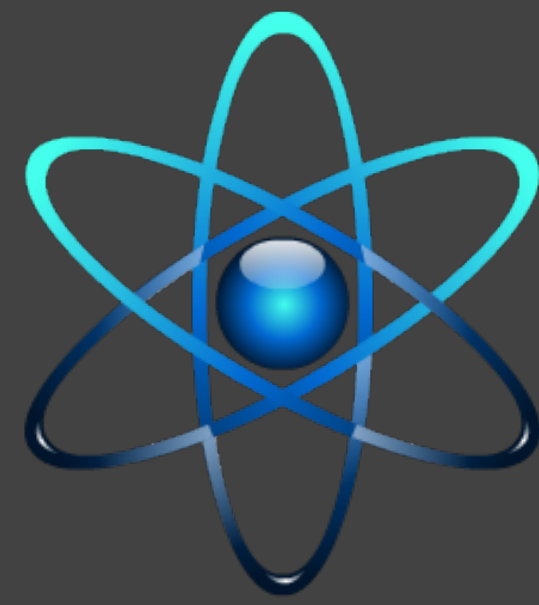
## Sarcasm generation

Chakrabarty et al. 2020 @ ACL 2020

## Simile generation

Chakrabarty et al. 2020 @ EMNLP 2020





# ATOMIC: An Atlas of Machine Commonsense for If-Then Reasoning

Maarten Sap

AAAI 2019



Ronan LeBras



Emily Allaway



Chandra Bhagavatula



Nicholas Lourie



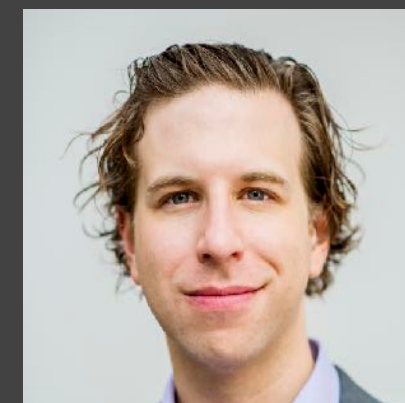
Hannah Rashkin



Brendan Roof



Noah Smith

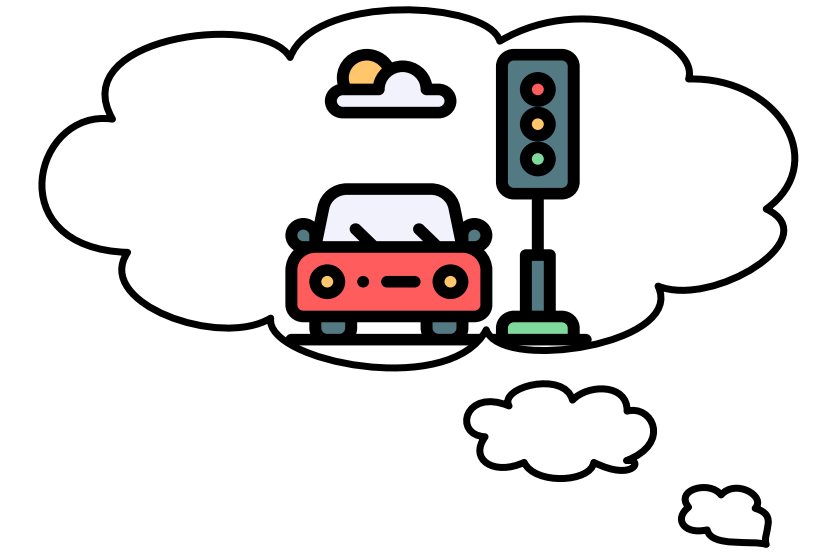


Me





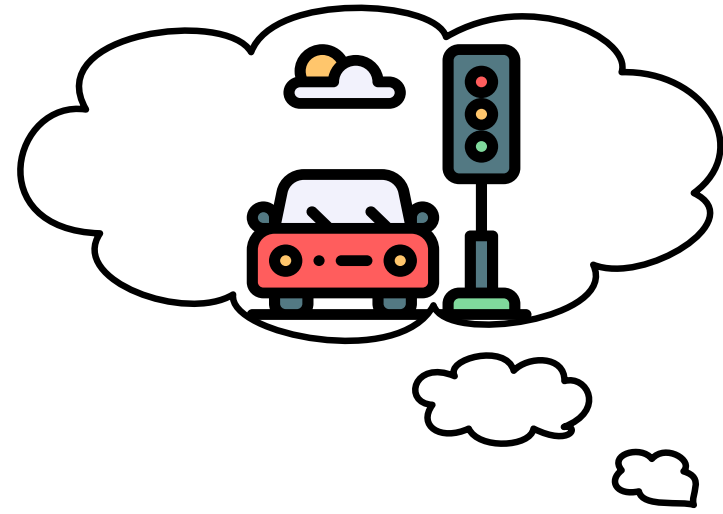
# Definition of Common Sense



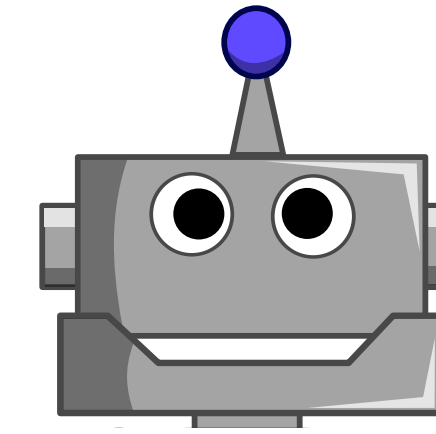
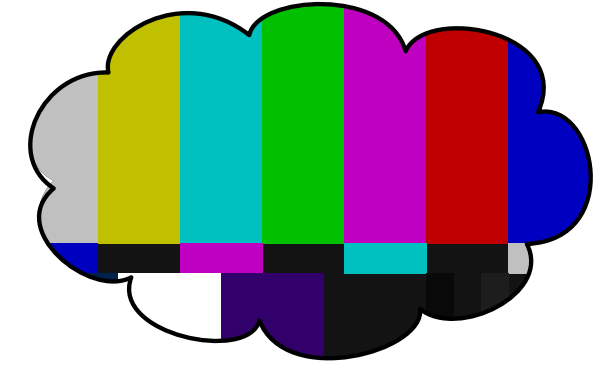
- the basic level of **practical knowledge** and **reasoning**
- concerning **everyday situations** and **events**
- that are **commonly** shared among **most** people.



For example, it's ok to keep the closet door open,  
but it's not ok to keep the fridge door open,  
as the food inside might go bad.



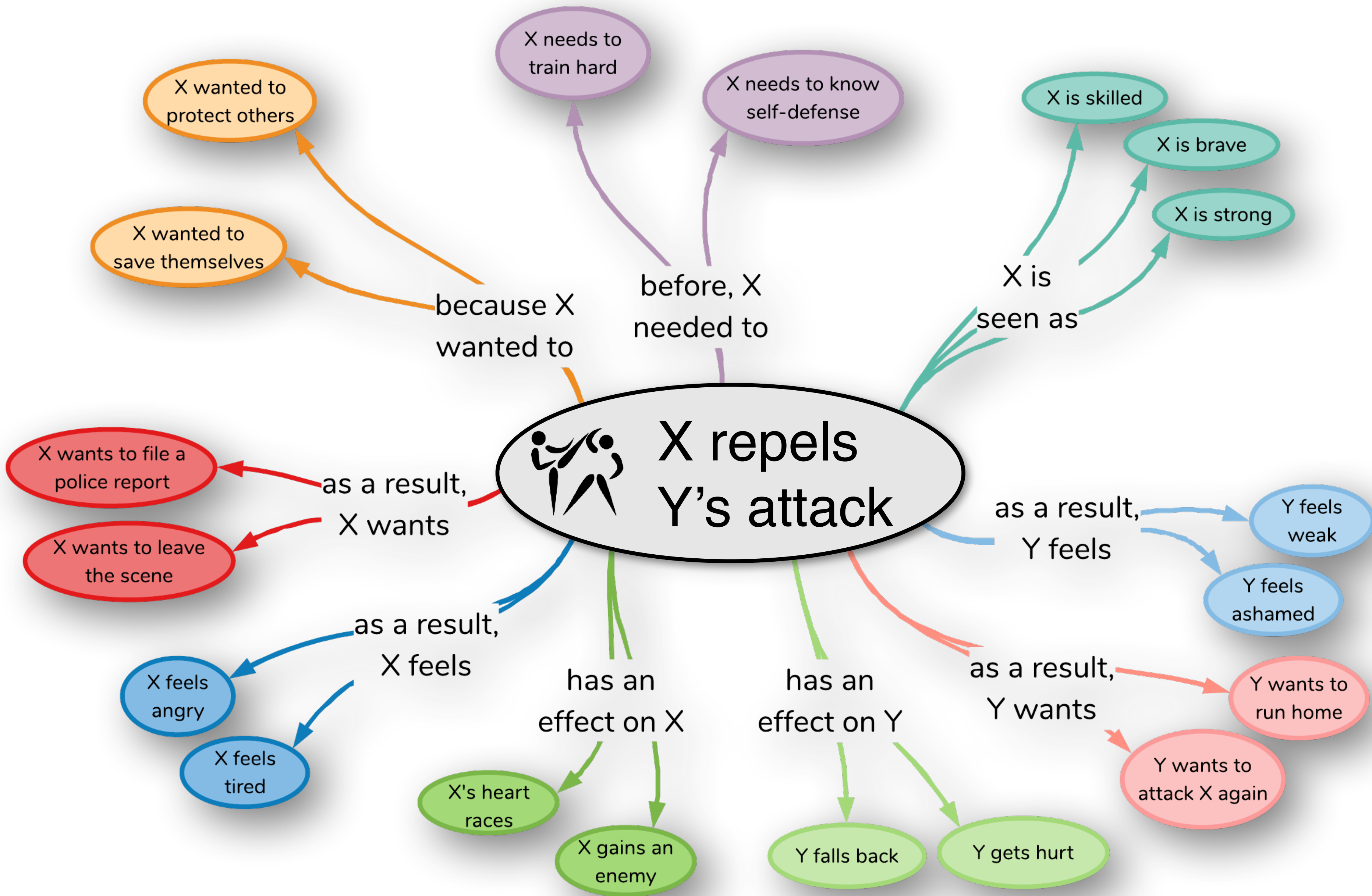
Essential for humans to live and interact with each other in a reasonable and safe way.



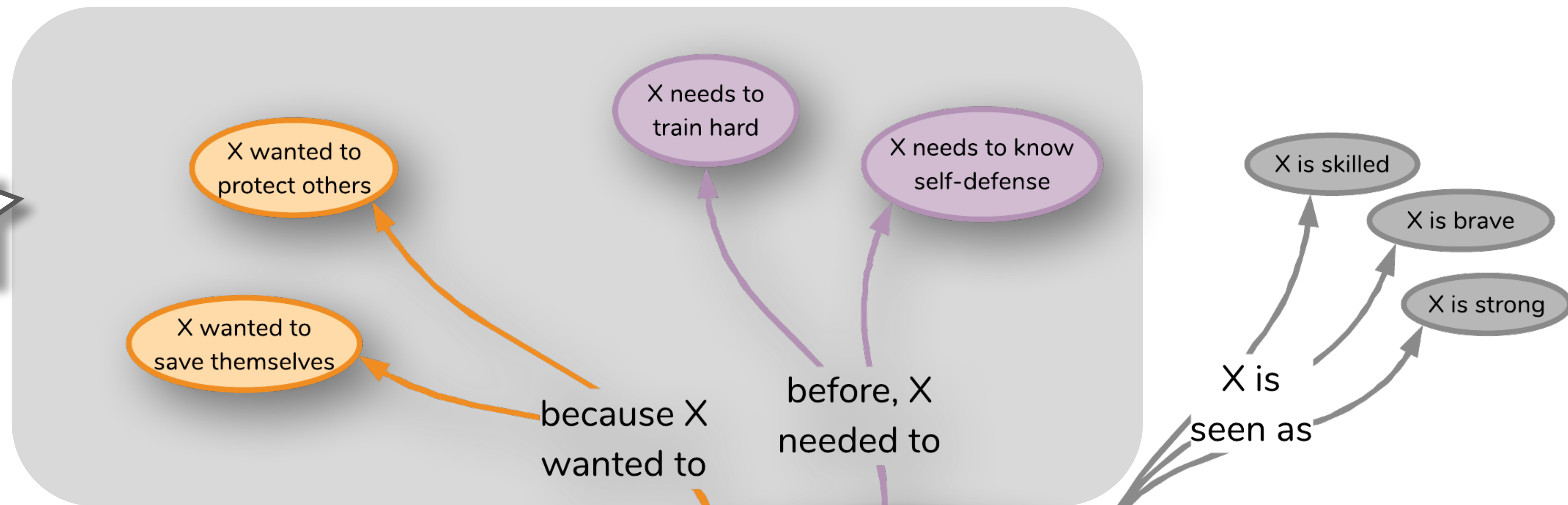
Essential for AI to understand human needs and actions better

For example, it's ok to keep the closet door open,  
but it's not ok to keep the fridge door open,  
as the food inside might go bad.

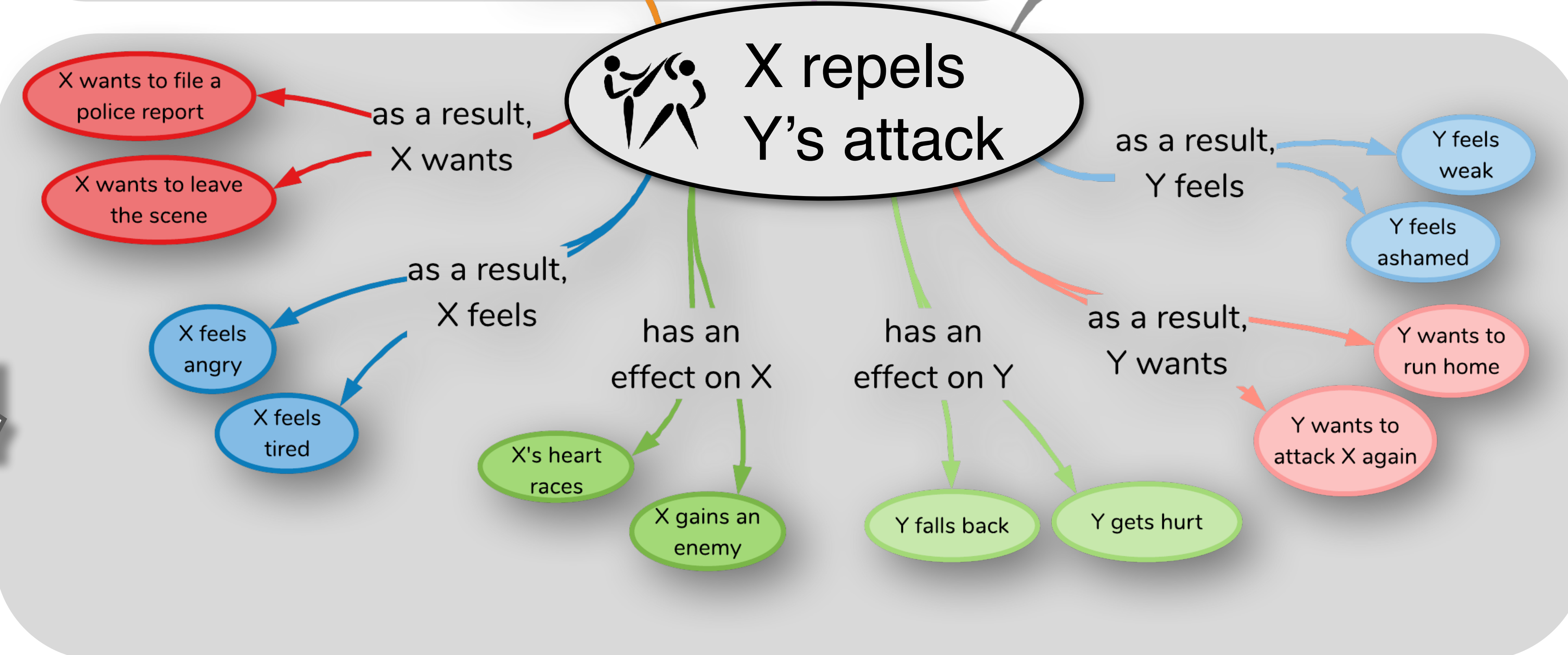




## Causes



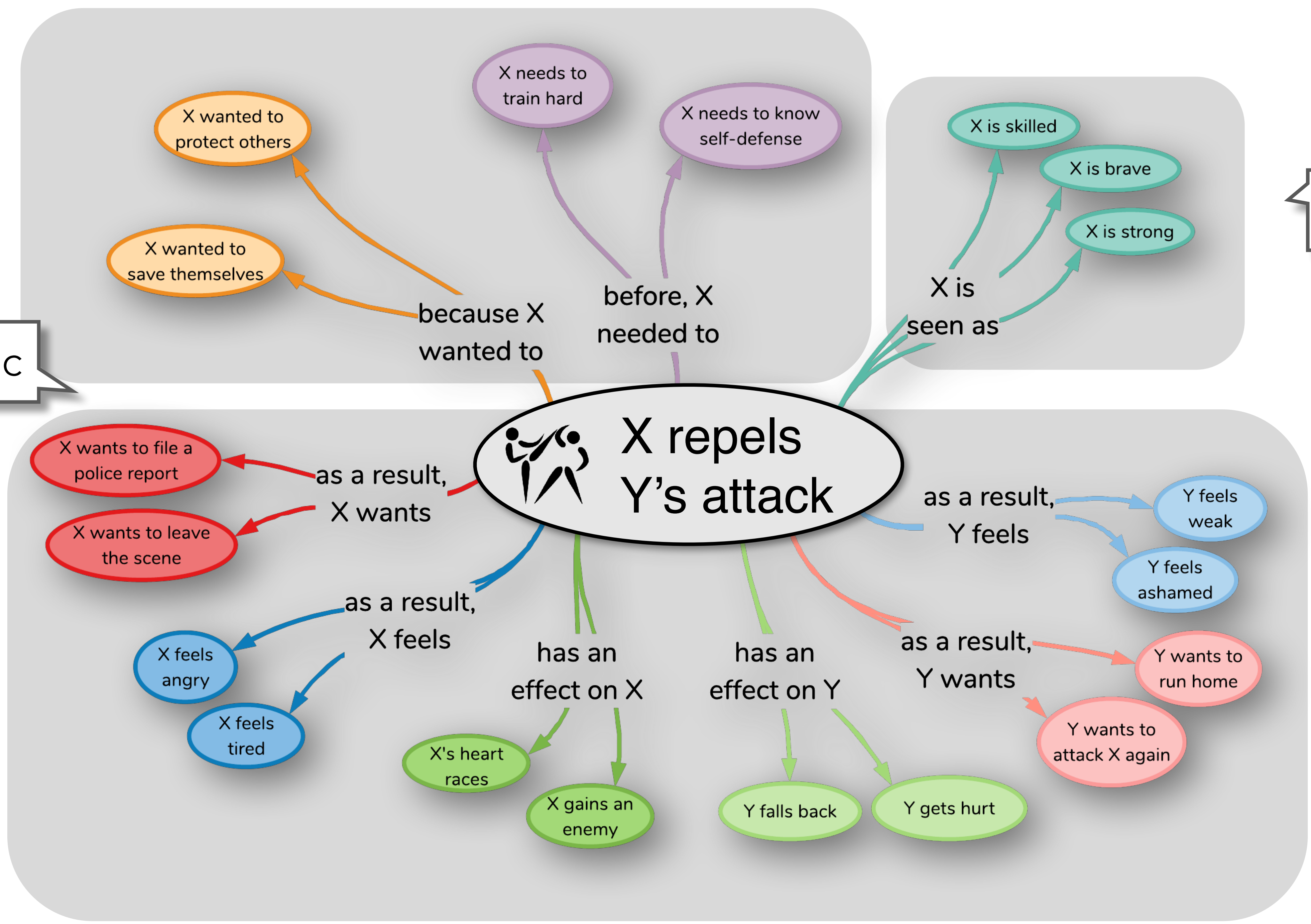
## Effects





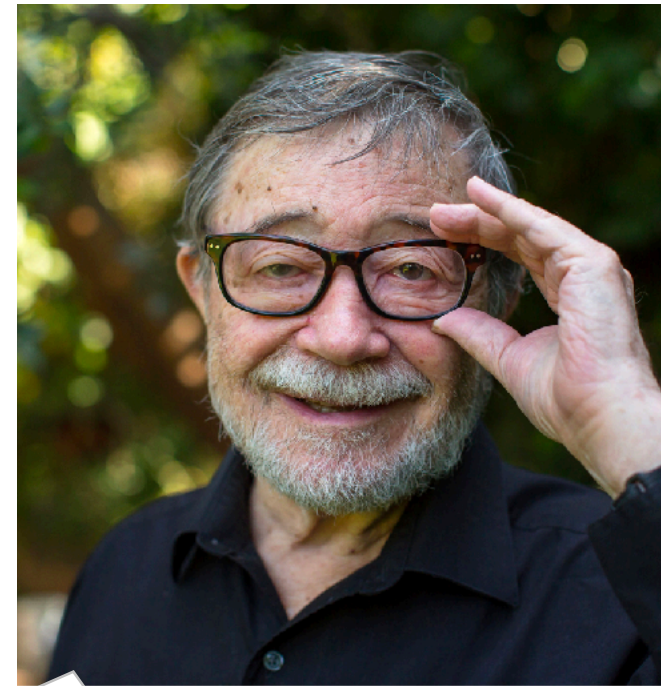
Dynamic

Static





# "Cause and Effect"




"Deep learning, I see they're all stuck there on the level of associations. Curve fitting."

"To build truly intelligent machines, teach them cause and effect"

Copyrighted Material

JUDEA PEARL  
*WINNER OF THE TURING AWARD*  
AND DANA MACKENZIE

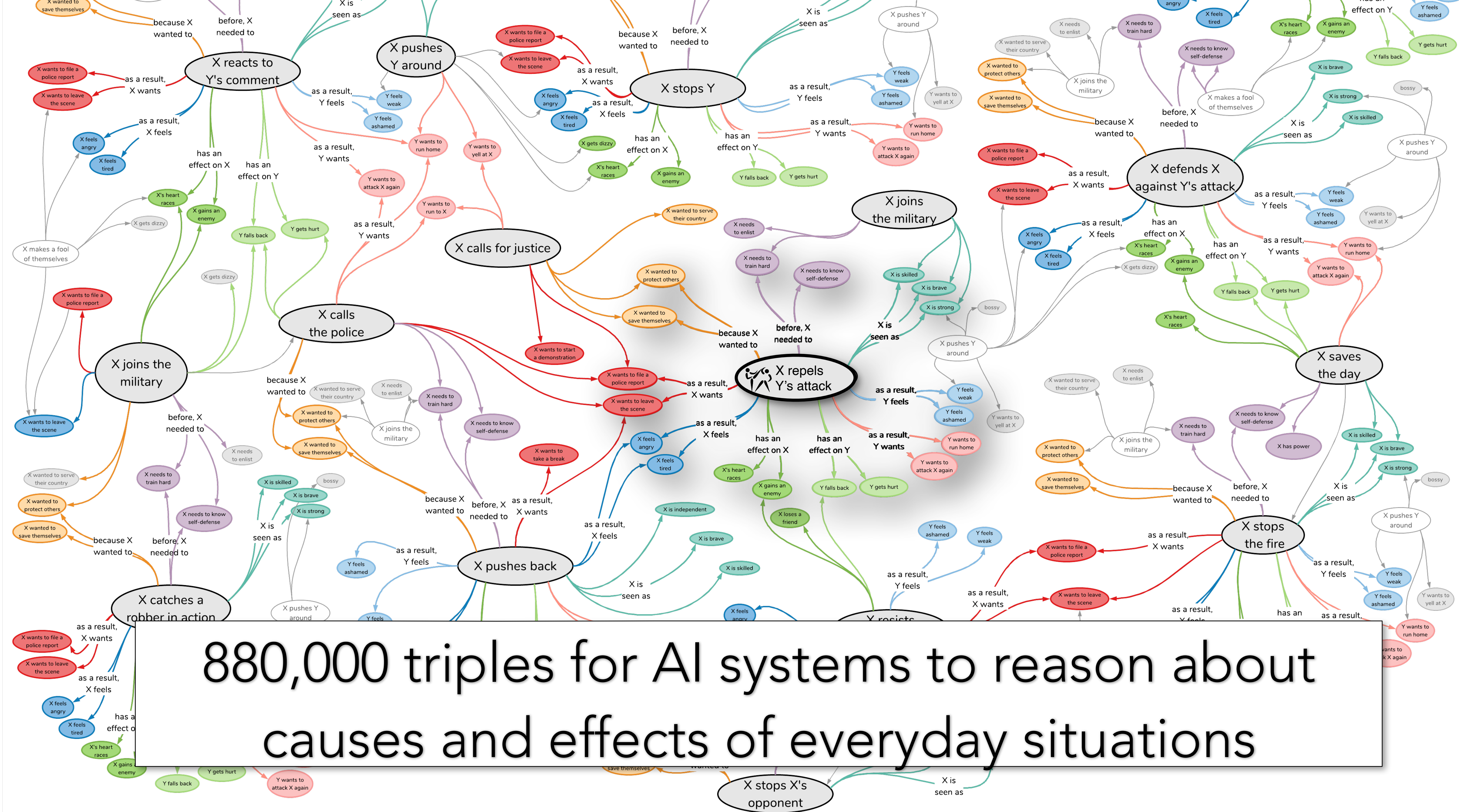
THE  
BOOK OF  
WHY

$\alpha$    $\beta$

THE NEW SCIENCE  
OF CAUSE AND EFFECT

Copyrighted Material







# How to acquire causes and effects at scale?

From unlabeled text?



- Reporting bias (Gordon & Van Durme, 2013)



*Murdering is 4x as  
common as exhaling?*



- Commonsense knowledge is not often written (Grice, 1975)



**ATOMIC**: Crowdsourced commonsense knowledge around **event prompts** using **natural language**

Browse ATOMIC:  
**<https://tinyurl.com/atomic-commonsense>**



# Existing commonsense knowledge bases

Knowledge of "what"

(taxonomic: **A isA B**; Davis and Marcus, 2015)

Represented in logical forms

(except ConceptNet)

```
event := (forall (e) (iff (event e) (or (exists (e1 e2) (and
(nequal e1 e2) (change' e e1 e2))) (exists (e1) (subevent
e1 e))))))
```

99% taxonomic

OpenCyc  
(Lenat, 1995)

EventNet  
(Espinosa &  
Lieberman,  
2005)

98% taxonomic

ConceptNet  
(Liu & Singh, 2004)

Formal Theory of  
Commonsense Psychology  
(Gordon & Hobbs, 2017)

Knowledge of "why" and "how"

(inferential: causes and effects)

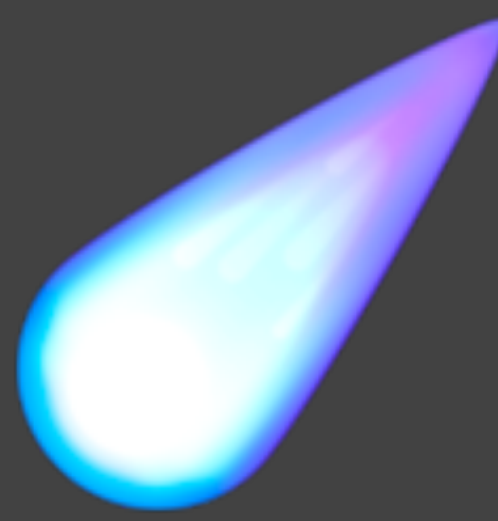
Represented in natural language

(how humans talk and think)

Visual Commonsense  
Graphs (Park et al., 2020)

ATOMIC  
(Sap et al.; 2019)





# COMeT: Commonsense Transformers for Automatic Knowledge Graph Construction

ACL 2019

Antoine  
Bosselut



Hannah  
Rashkin



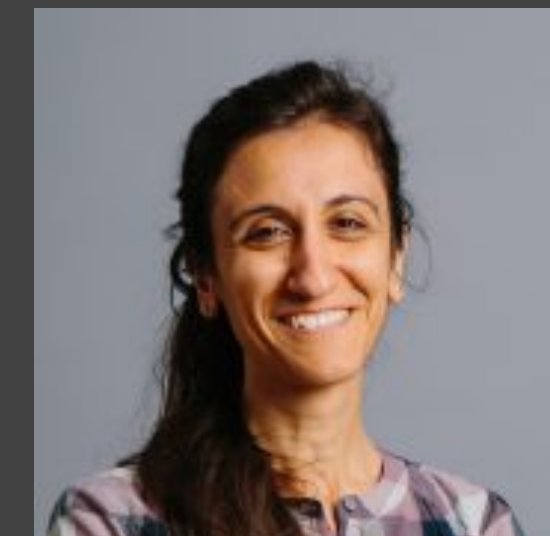
Maarten  
Sap



Chaitanya  
Malaviya



Asli  
Çelikyilmaz



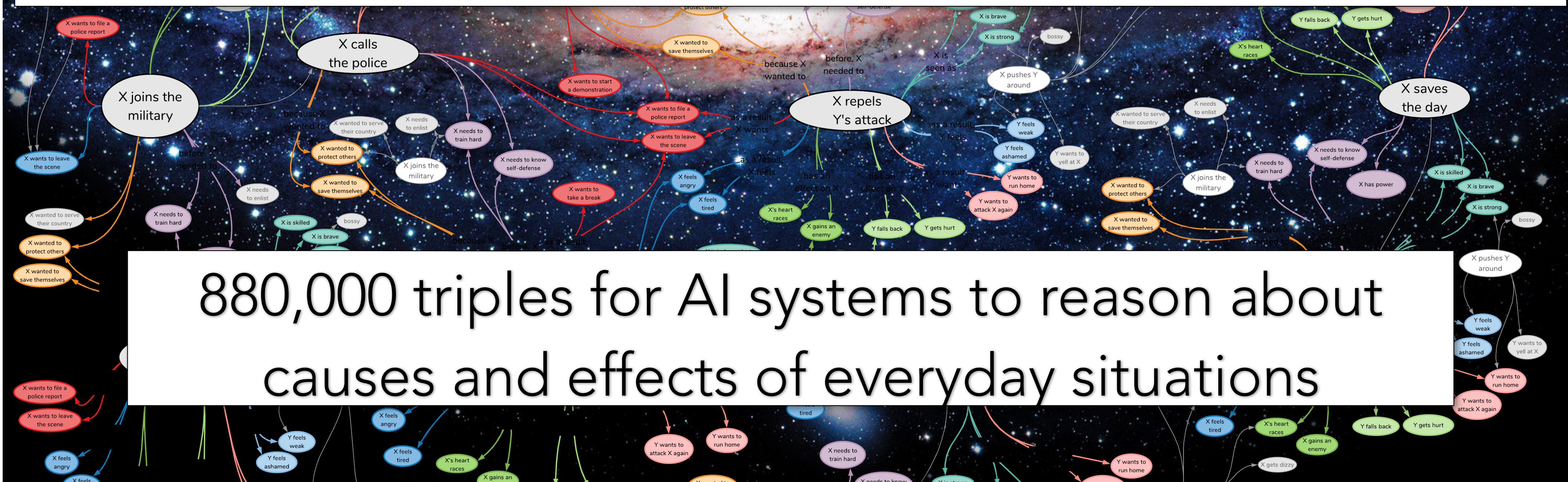
Me



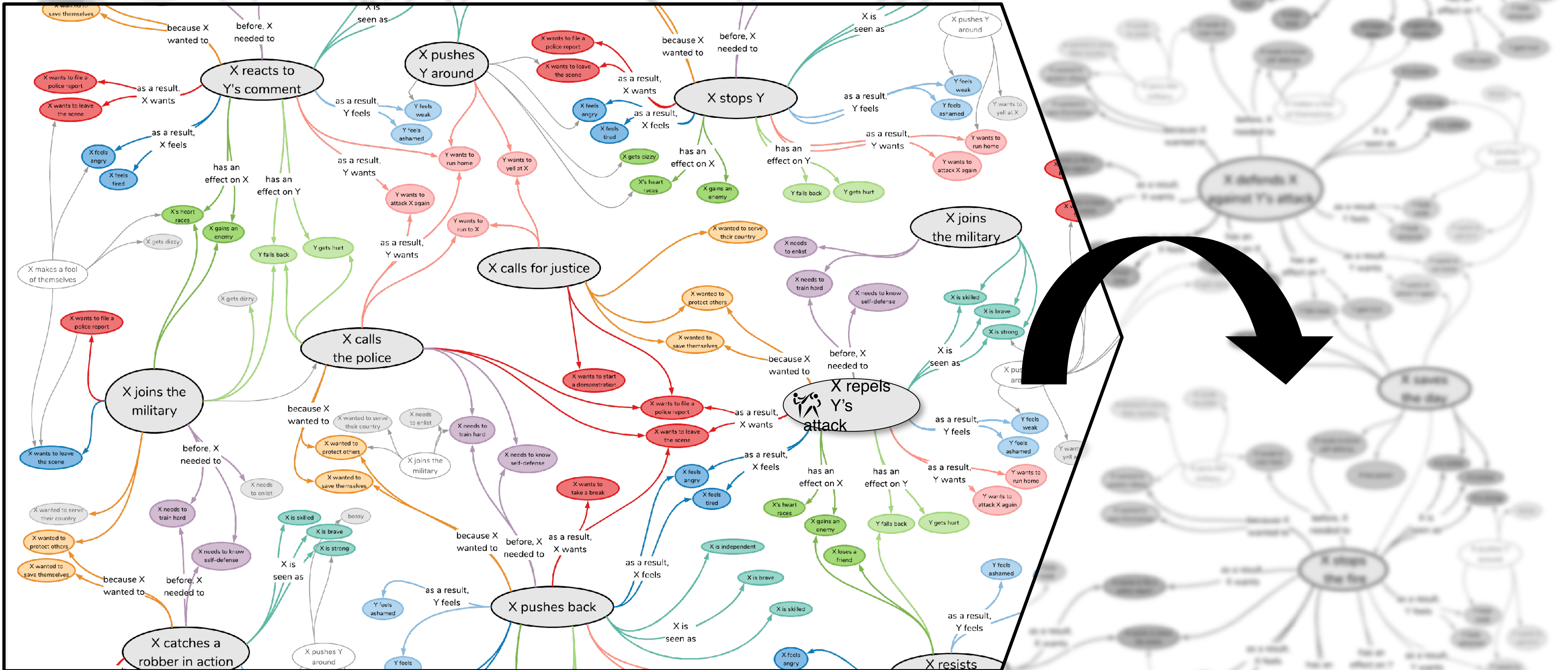


Can we reason about commonsense knowledge without storing all of them explicitly?

Transfer learning from language (self-supervised) to knowledge (supervised)?








Goal: to teach models to reason about causes and effects of new ATOMIC events

# COMeT

Transfer learning from  
language to knowledge!



*Reporting bias issues, but  
provides signal about  
which events are similar*



# COMeT

**PersonX sails across the Atlantic**



***context event***

**<xNeed>**



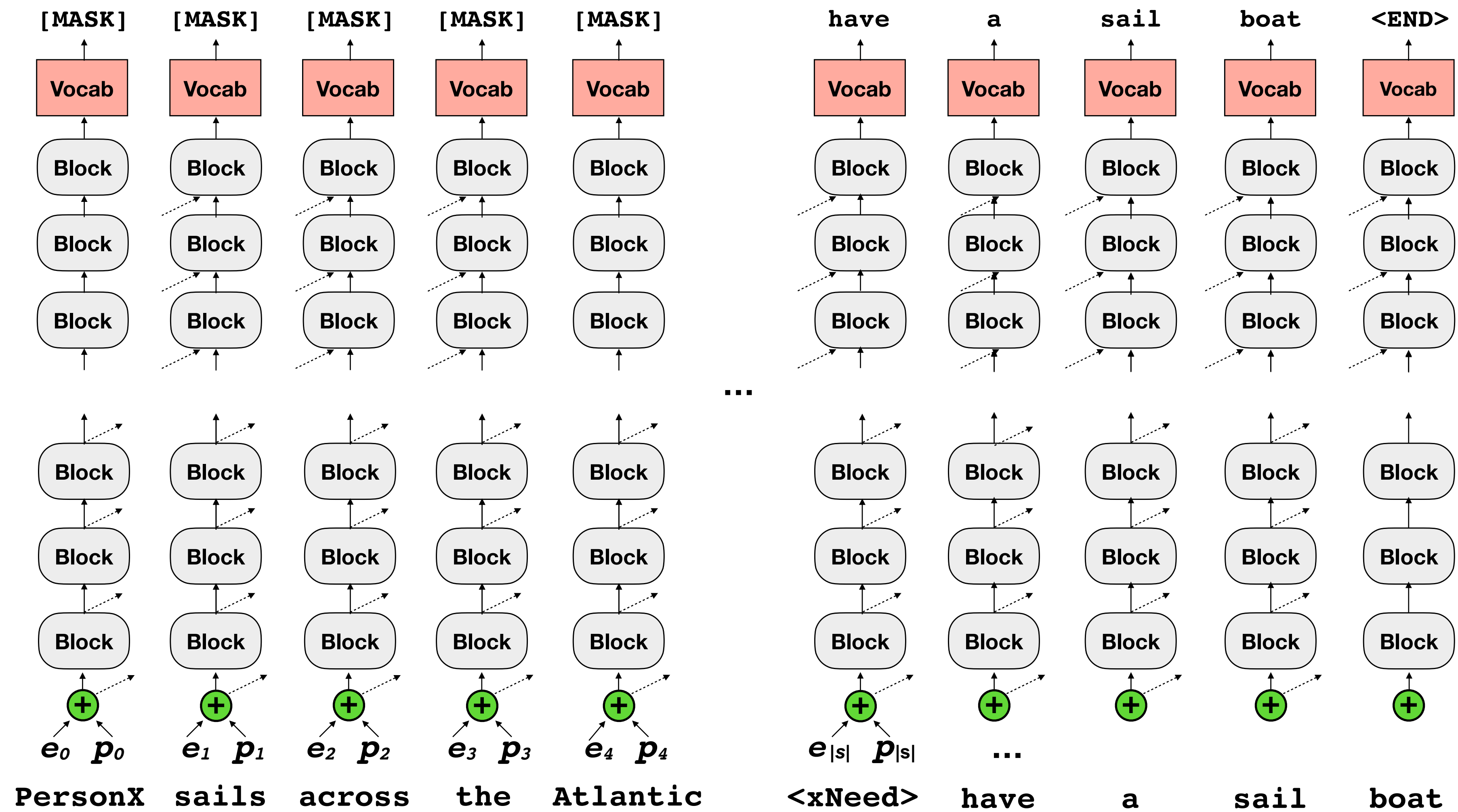
***desired  
common  
sense  
inference***

**have a sailboat**



***commonsense  
prediction***

# COMeT





Turns out, COMET generalizes well  
on out-of-domain examples

(which we realized only after publishing our ACL 2019 paper... )

# Demo: <https://mosaickg.apps.allenai.org/>

**AI2** Allen Institute for AI



## Mosaic Knowledge Graphs

COMeT ^

Events (ATOMIC)

Concepts (ConceptNet)

Images (VisualCOMET)

Knowledge Graph ^

ATOMIC

About

### Commonsense Inferences about Events (COMmonsenseE Transformers on ATOMIC)

A knowledge base construction e  
graphs. By training on a seed set

Write a short text to reason about

### Explore

sarah repel's jack's attack in a chess game

Predict

Try: [PersonX acts quickly](#), [PersonX is a big deal](#), [My boss is very good](#)

### COMeT Predictions Graph

The model has predicted these relationships for 'sarah repel's jack's attack in a chess game'  
[View 'sarah repel's jack's attack in a chess game' in the ATOMIC dataset](#)

Causes for PersonX

Because PersonX **wanted**

- to win
- to win the game
- to play a game
- to win a game
- none

Click here for COMET demo  
(COMET trained on ATOMIC)



# Demo: <https://mosaickg.apps.allenai.org/>

AI2 Allen Institute for AI



## Mosaic Knowledge Graphs

COMeT

Events (ATOMIC)

Concepts (ConceptNet)

Images (VisualCOMET)

Knowledge Graph

ATOMIC

About

### Commonsense Inferences about Events (COMmonsenseE Transformers on ATOMIC)

A knowledge base construction e  
graphs. By training on a seed set

Write a short text to reason about

### Explore

sarah repel's jack's attack in a chess game

Predict

Try: PersonX acts quickly, PersonX is a big deal, My boss is very good

### COMeT Predictions Graph

The model has predicted these relationships for 'sarah repel's jack's attack in a chess game'  
[View 'sarah repel's jack's attack in a chess game' in the ATOMIC dataset](#)

Causes for PersonX

Because PersonX **wanted**

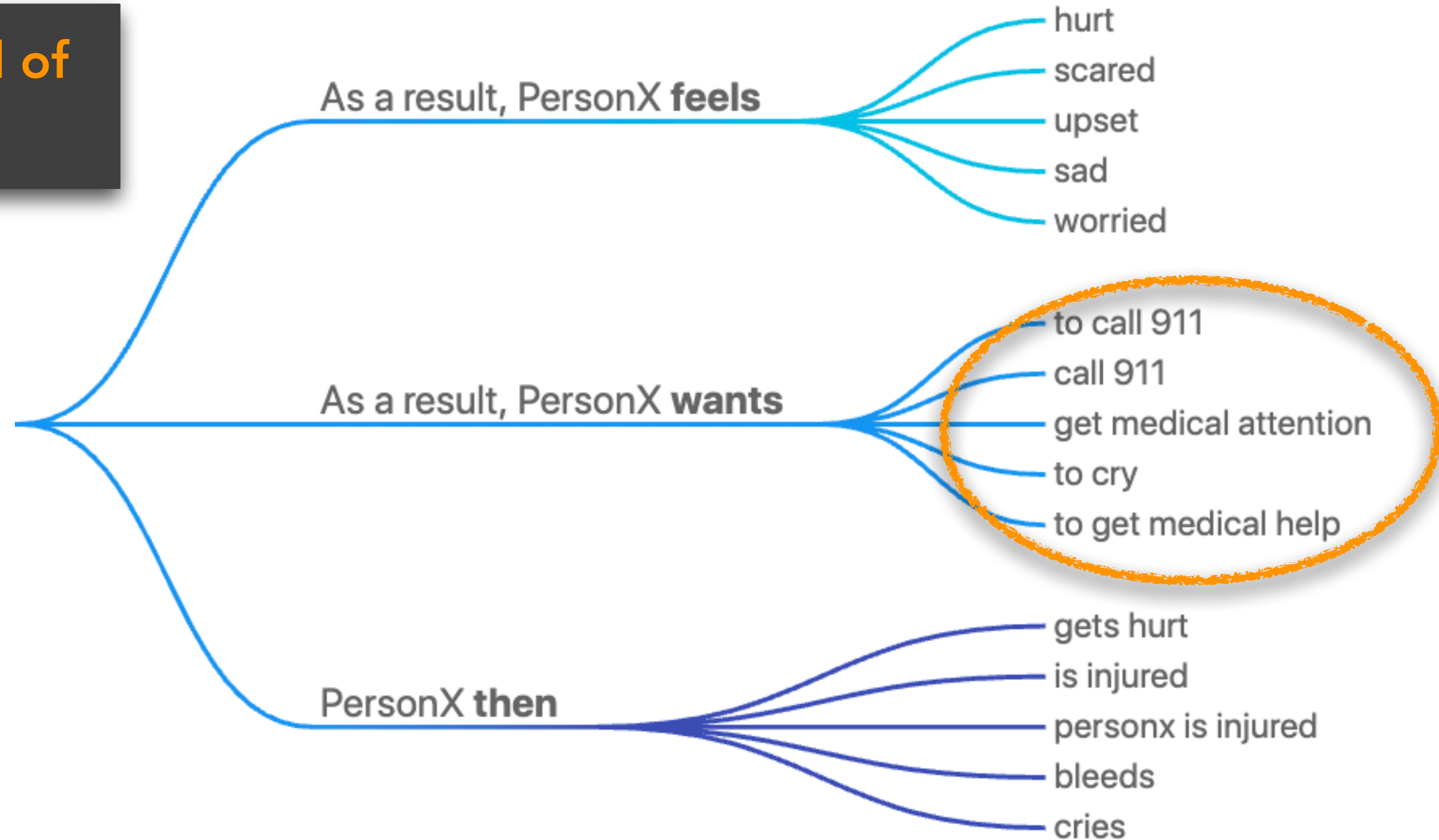
- to win
- to win the game
- to play a game
- to win a game
- none

Click here for COMeT demo  
(COMeT trained on ATOMIC)

Click here for COMeT demo  
(COMeT trained on ConceptNet)

# John gets into an accident

"John" instead of  
"Person X"





(Some amount of)

# Abstraction by COMET?

# Demo: <https://mosaickg.apps.allenai.org/>

**A12** Allen Institute for AI



## Mosaic Knowledge Graphs

COMeT ^

Events (ATOMIC)

Concepts (ConceptNet)

Images (VisualCOMET)

Knowledge Graph ^

ATOMIC

About

### Explore the ATOMIC Knowledge Graph

An atlas of everyday commonsense reasoning, organized through 877k textual descriptions of inferential knowledge. Compared to existing resources... [more](#) | [read the paper](#) | [download the data](#)

### Explore

PersonX is self conscious

trn

PersonX is very self conscious

trn

PersonX is being controlled by PersonY's subconscious

trn

PersonX is controlled by PersonY's subconscious

trn

PersonX is knocked unconscious

trn

Select an event to see relationships in the dataset

1. Type any one content word ("conscious"), then the demo lists all events with the query word in the KG
2. Choose one of the events in the dropdown, hit enter

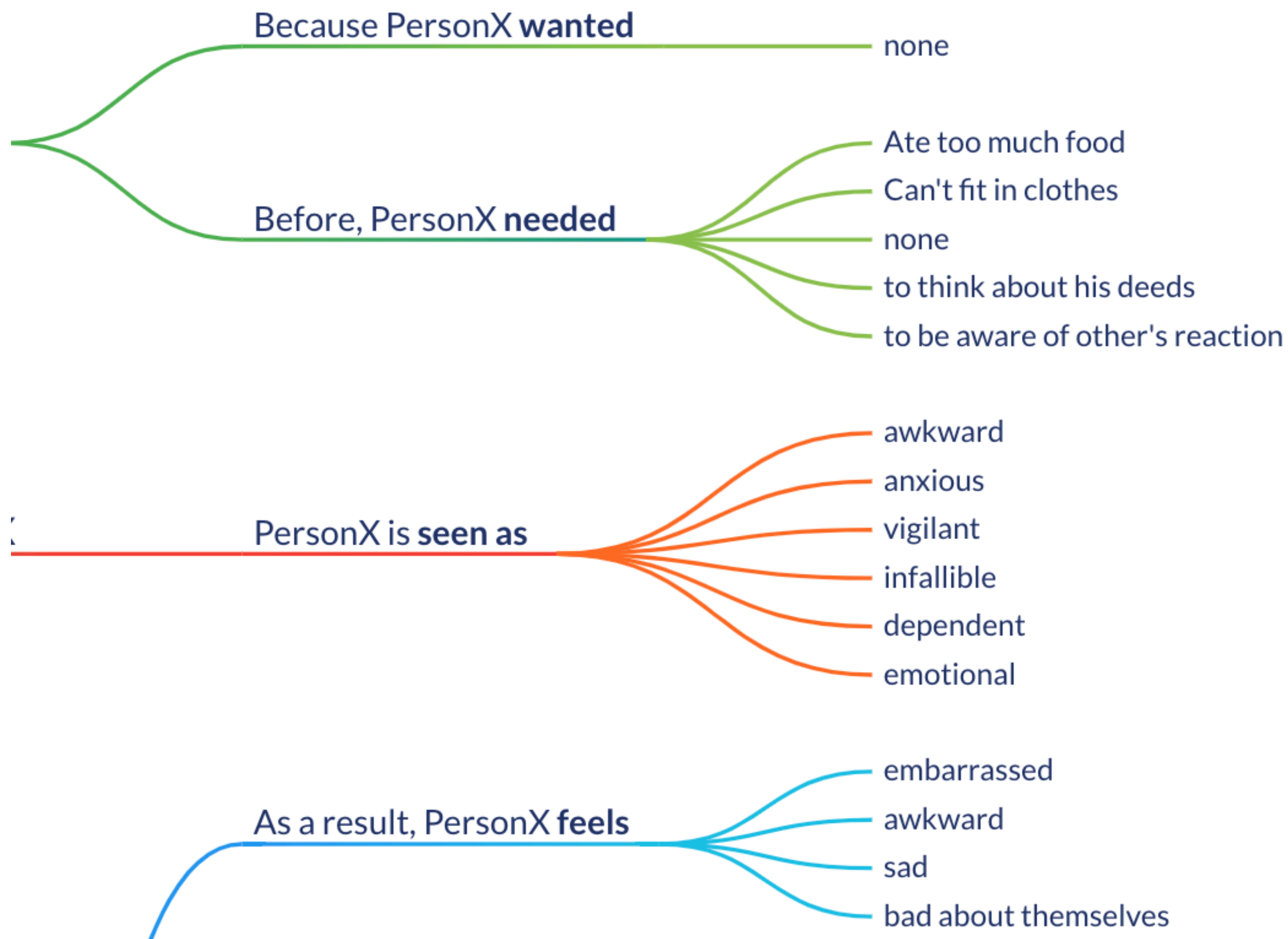
Click here to search ATOMIC knowledge graph





ATOMIC

# "PersonX is self conscious"





ATOMIC

# "PersonX is self conscious"



COMET

Because PersonX wanted

none

Before, PersonX needed

Ate too much food

Can't fit in clothes

none

to think about his deeds

to be aware of other's reaction

PersonX is seen as

awkward

anxious

vigilant

infallible

dependent

emotional

As a result, PersonX feels

embarrassed

awkward

sad

bad about themselves

|

to look good

to show off

to be self - conscious

|

none

to do something embarrassing

to do something wrong

to be nervous

to have a lot of work

|

nervous

shy

insecure

timid

awkward

|

nervous

embarrassed

awkward

ashamed

uncomfortable





ATOMIC

"PersonX is self conscious"



COMET

PersonX wanted

none

to look good

to show off

to be self - conscious

PersonX needed

Ate too much food

Can't fit in clothes

none

to think about his deeds

to be aware of other's reaction

none

to do something embarrassing

to do something wrong

to be nervous

to have a lot of work

PersonX is seen as

awkward

anxious

vigilant

infallible

dependent

emotional

nervous

shy

insecure

timid

awkward

# "Sarah is conscious"

COMET

Recall that ATOMIC  
does NOT have any  
event about being  
**conscious**

ATOMIC has events  
only about being  
**Self-conscious**  
**Unconscious**

Because PersonX wanted

- none
- to be aware
- to observe something
- to observe
- to know what is going on

Before, PersonX needed

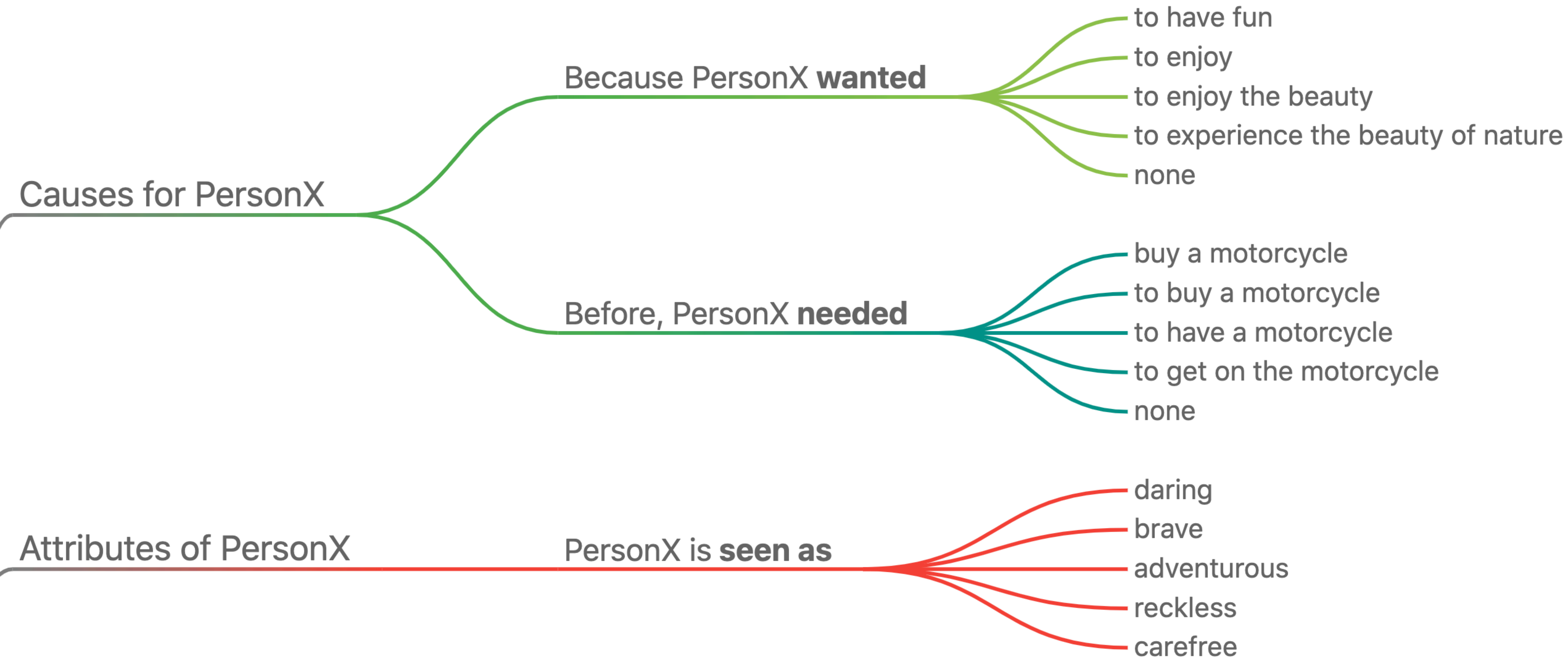
- none
- to be aware
- to hear something
- to have a headache
- to be aware of something

PersonX is seen as

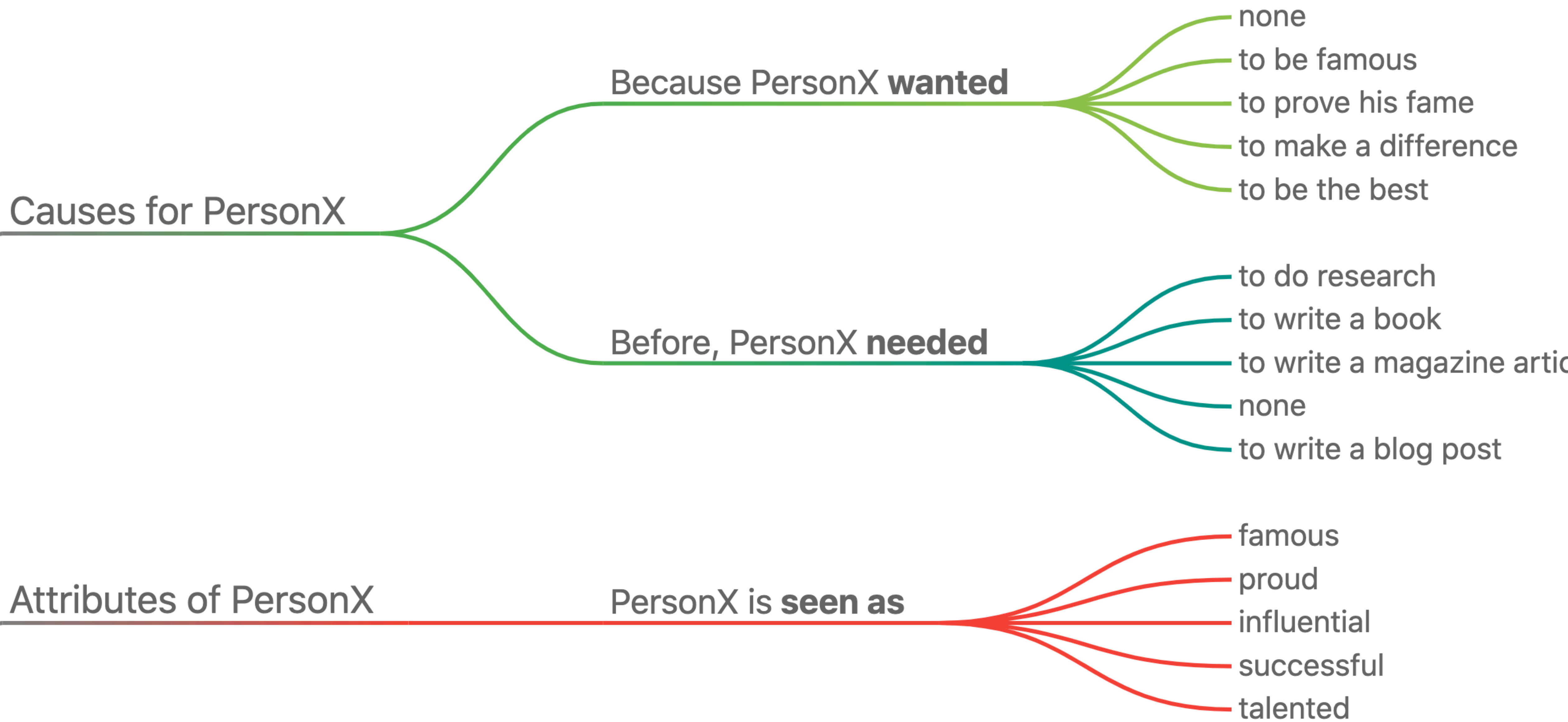
- aware
- conscious
- observant
- curious
- sensitive



# Sanja rides into the sunset on a motorcycle after solving AI.

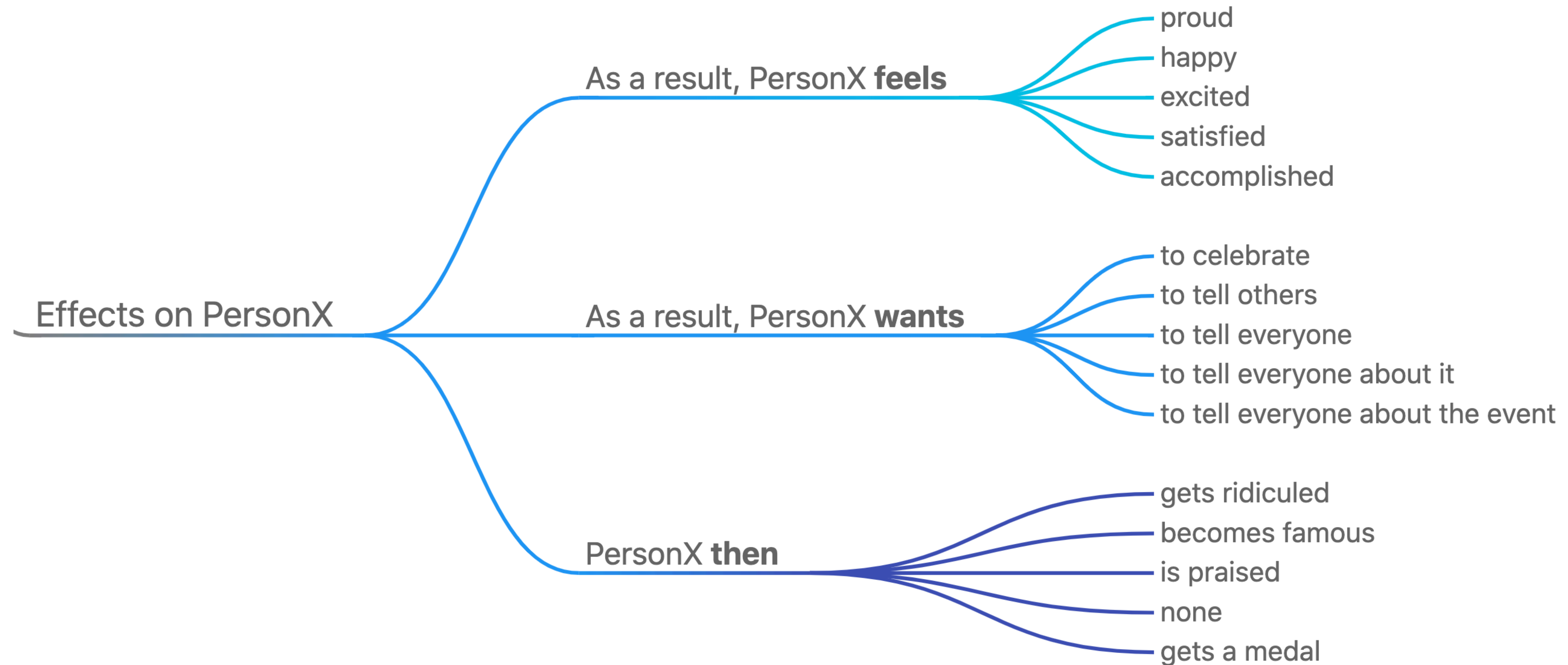


# Gary breaks the world record for most controversial tweet.





# Gary breaks the world record for most controversial tweet.





**Gary Marcus** @GaryMarcus · 1h



ps i realize there is a typo in my query. fixing the typo doesn't much help though:

## Completion

**what happens when you stack kindling and logs in a fireplace and then drop some matches is that you typically start a** ick. So, it's kind of ironic that the second day after my son was born, the fire in the living room had melted through the kindling. It's pretty neat."





**Yejin Choi**  
@YejinChoinka

Replying to @GaryMarcus

Gary, try [mosaickg.apps.allenai.org](https://mosaickg.apps.allenai.org) by typing "Gary stacks kindling and logs and drops some matches". Sorry I used deep learning... :)



focus on  
"causes and effects"  
(causal knowledge)

COMeT



NEURAL  
(generalizes well to  
compositional &  
unseen events)

(semi-) supervised learning of  
declarative knowledge

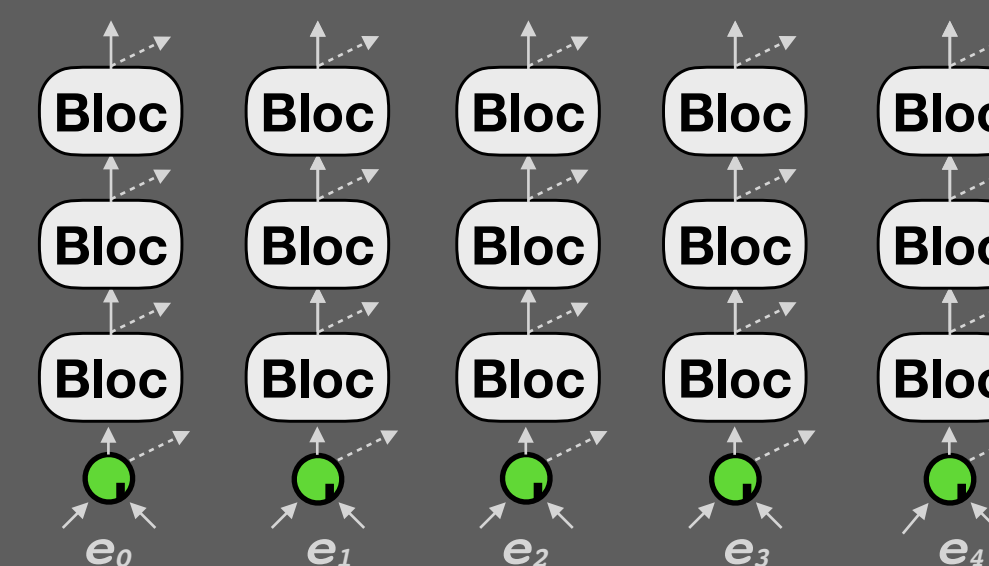
self supervised learning of  
observed knowledge

ATOMIC  
[tinyurl.com/atomic-](https://tinyurl.com/atomic-)



SYMBOLIC  
but in LANGUAGE  
(instead of LOGIC)

Language Models





# (COMET-) ATOMIC<sub>20</sub><sup>20</sup>:

## On Symbolic and Neural Commonsense Knowledge Graphs

— wait, doesn't GPT-3 know everything? —

To appear at AAAI 2021

Jena  
Hwang



Chandra  
Bhagavatula



Ronan  
Le Bras



Jeff  
Da



Keisuke  
Sakaguchi

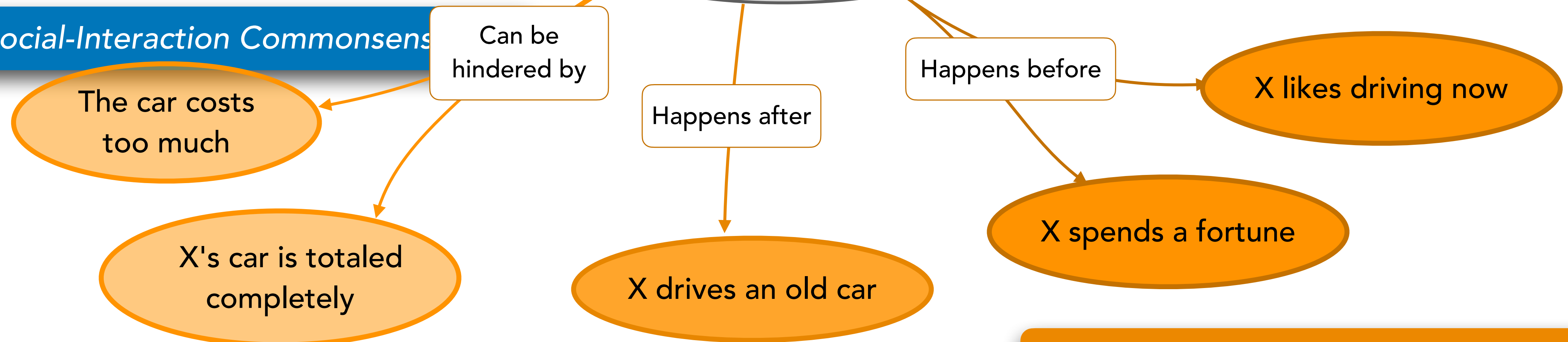
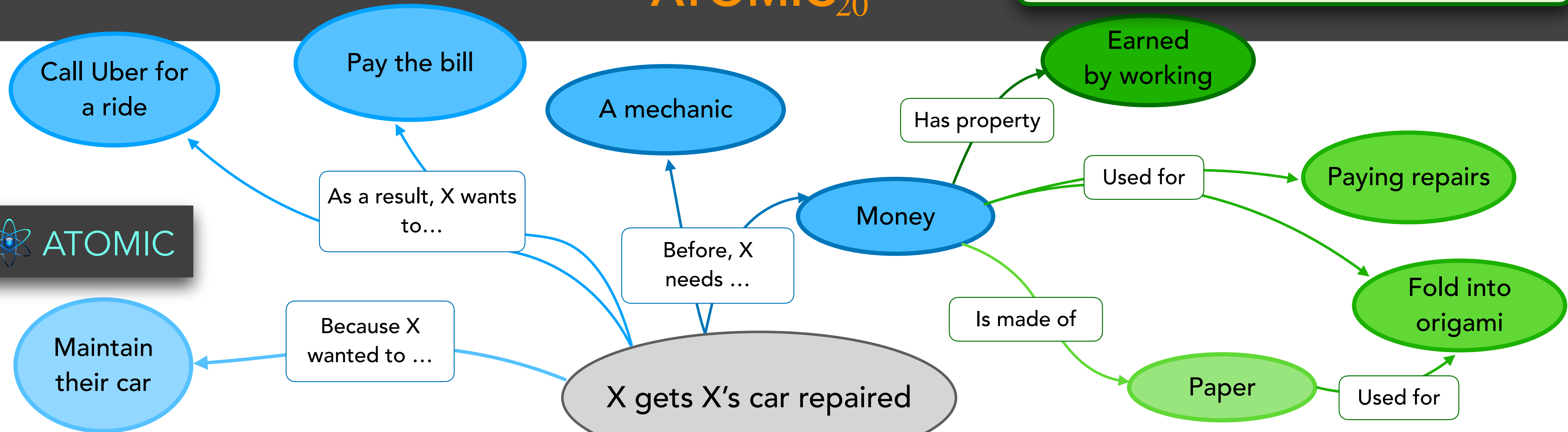


Antoine  
Bosseult



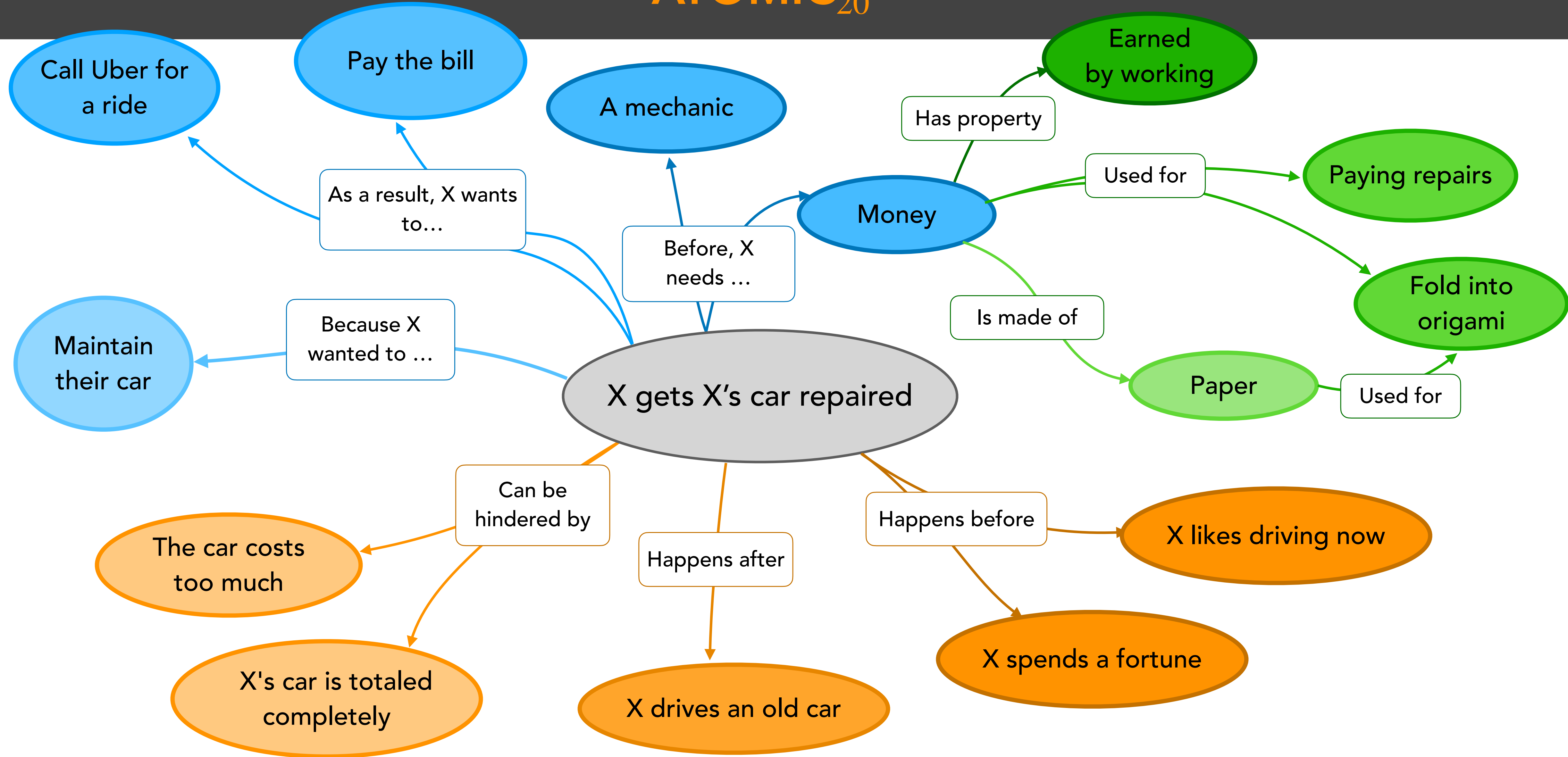
Me

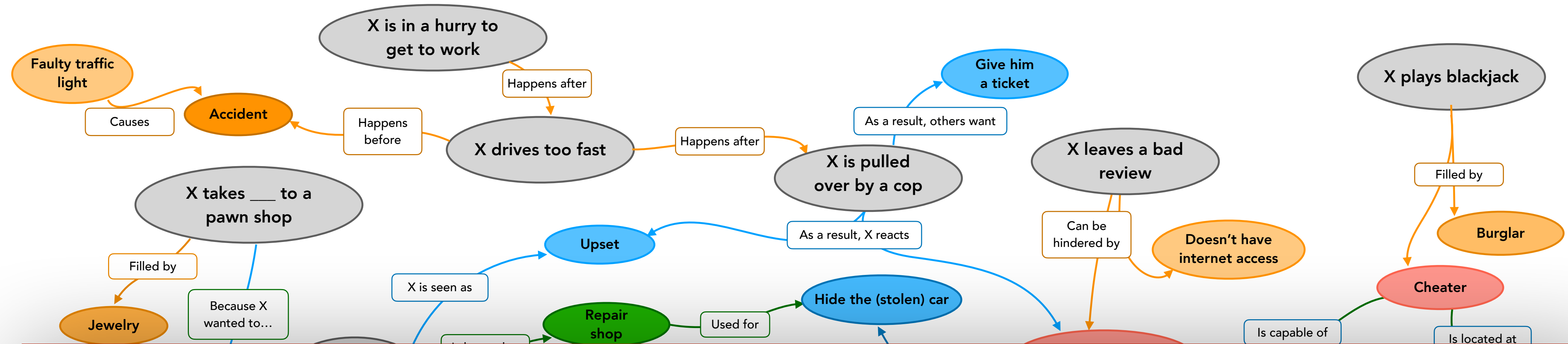




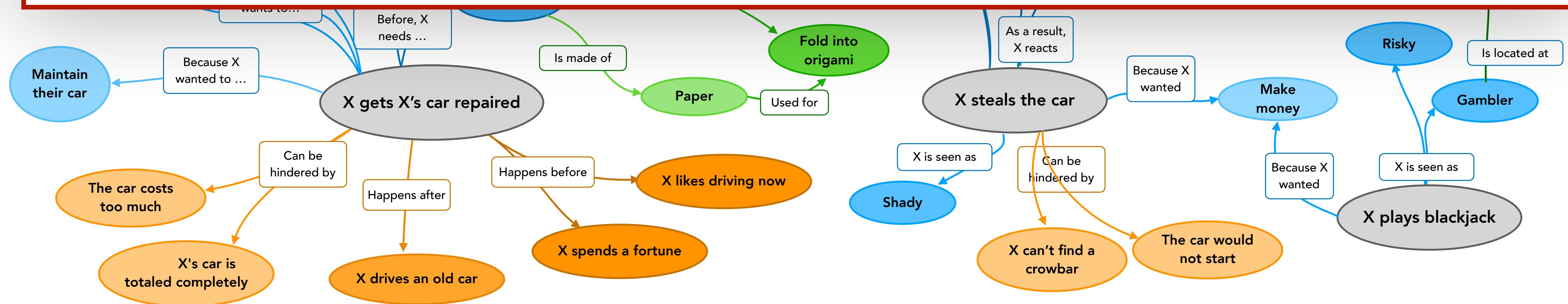


# ATOMIC<sup>20</sup><sub>20</sub>





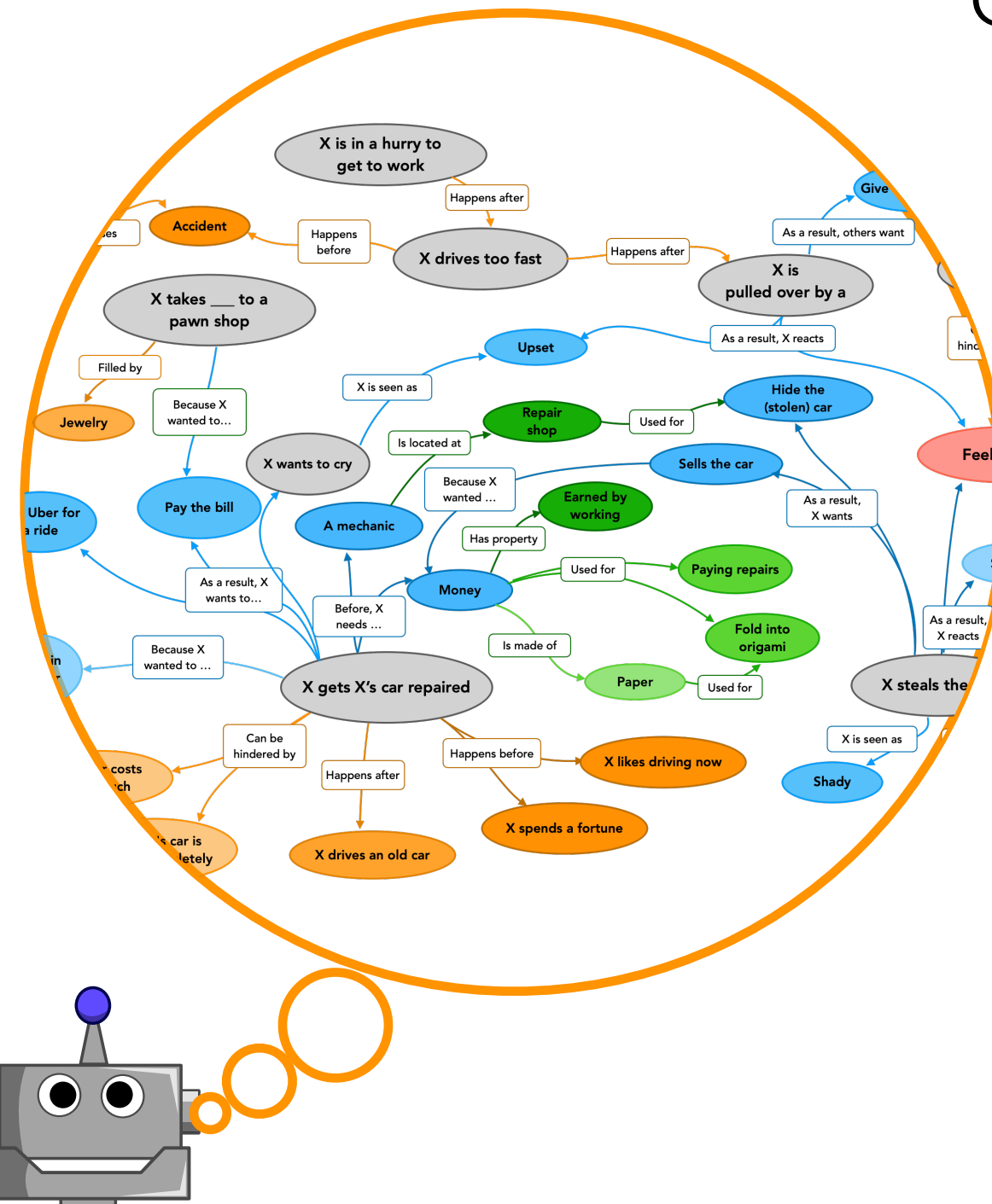
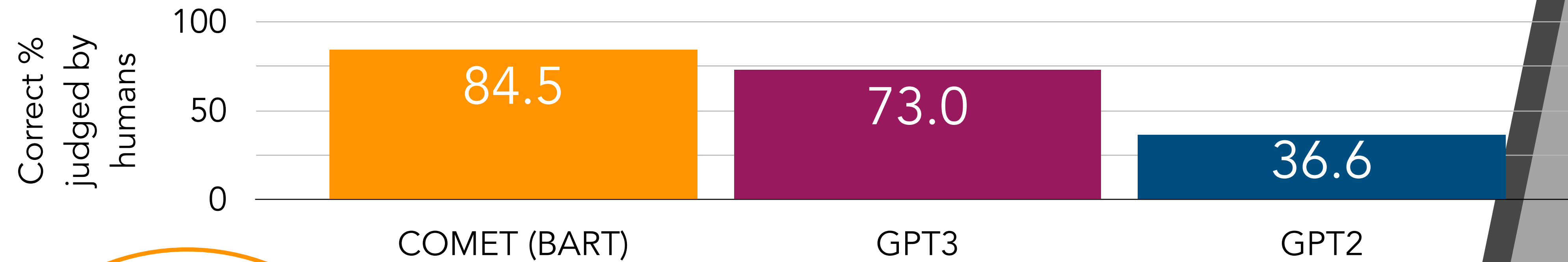
**1.33M commonsense if-then inferences**  
**23 relations (or inference types)**





## Knowledge Models

## Off-the-shelf Language Models



COMeT (BART): x435 smaller model (~400M parameters),  
informed by **ATOMIC**<sub>20</sub><sup>20</sup>

**GPT-3 (Few Shot): 175B parameters!!**  
pre-trained with a ton of web text (~500B tokens)

# Visual COMET:

Reasoning about the *Dynamic* Context of a *Still* Image

ECCV 2020

Jae Sung (James) Park



Chandra  
Bhagavatula



Roozbeh  
Mottaghi



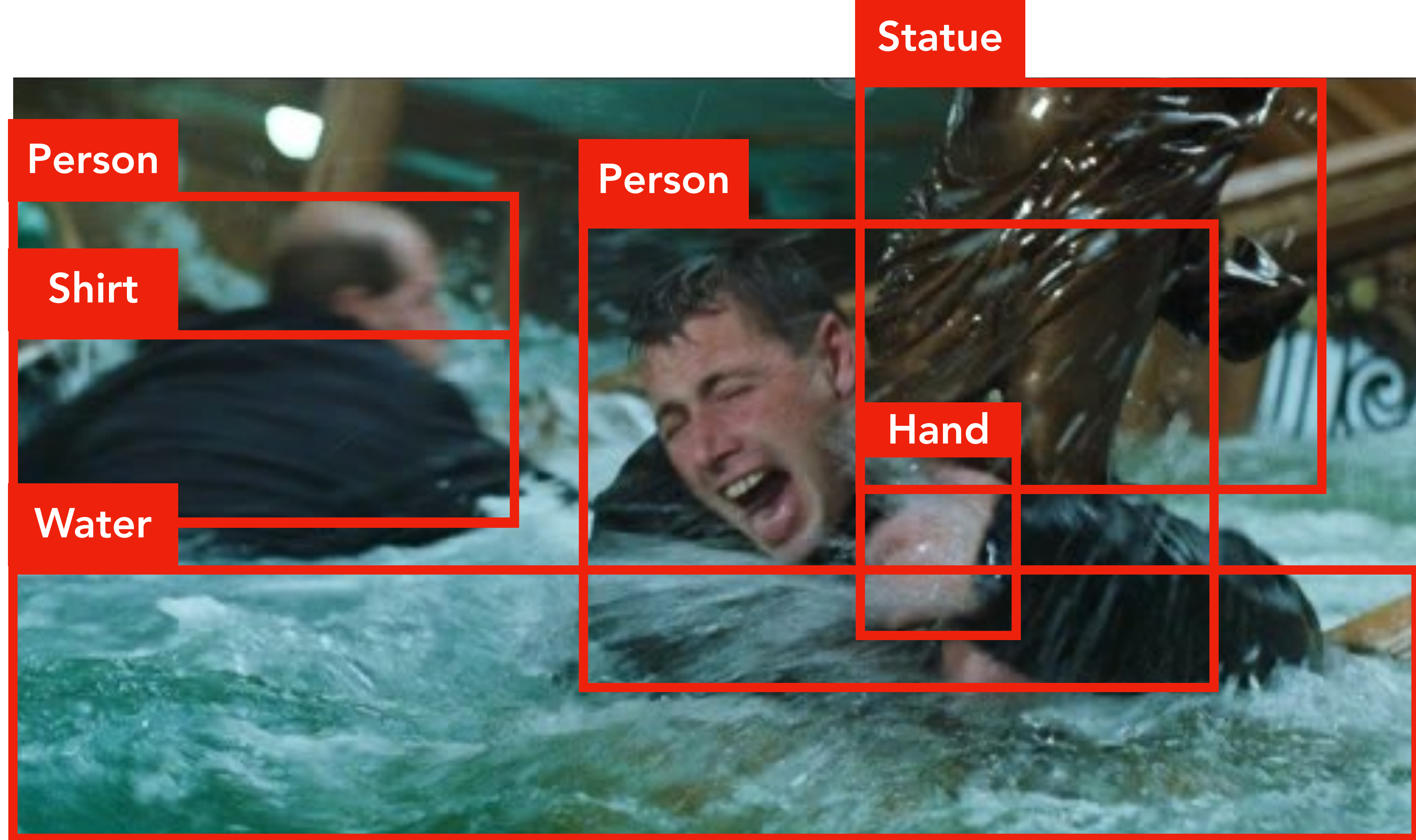
Ali  
Farhadi



Yejin  
Choi

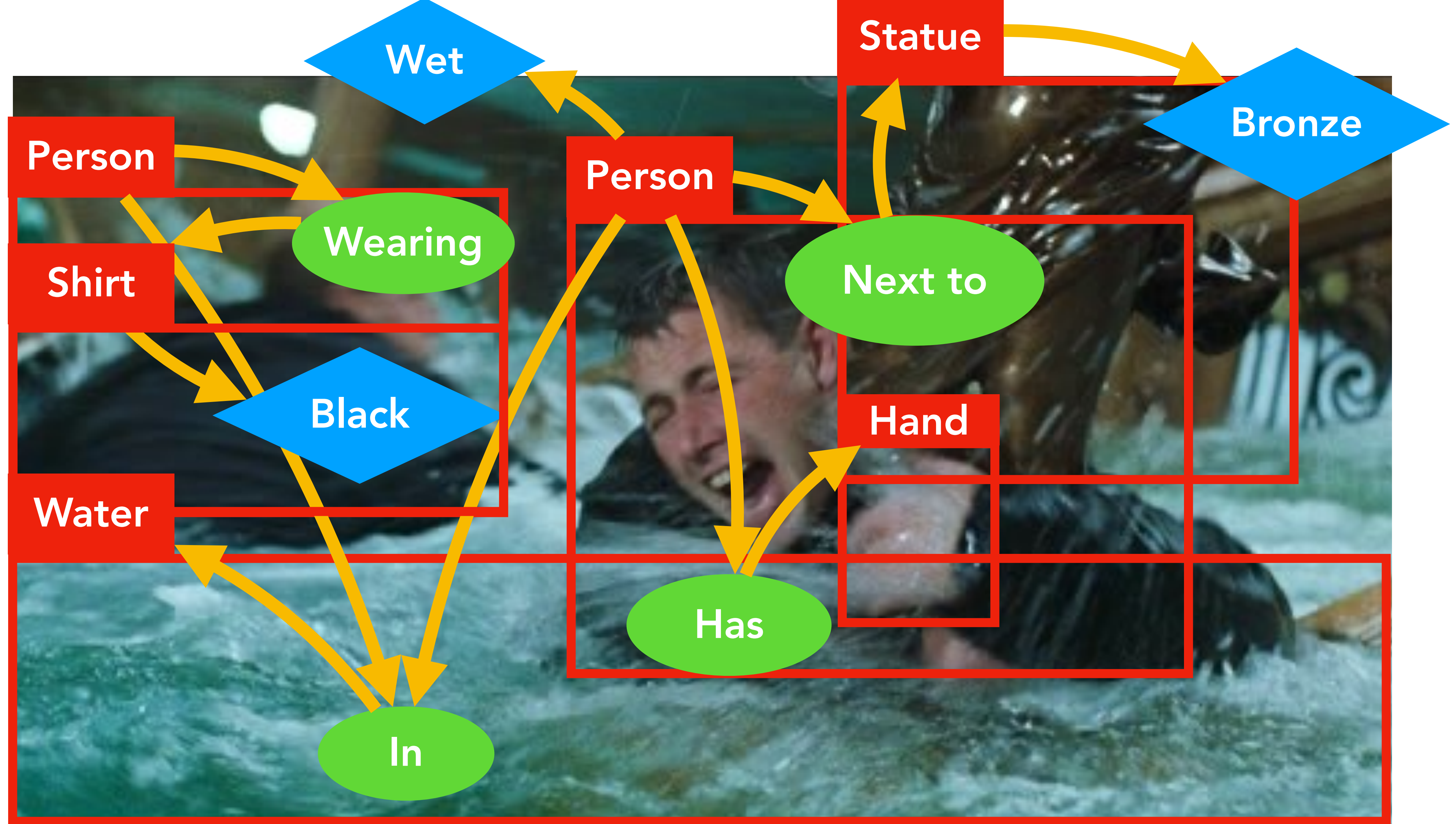






# Object Detection

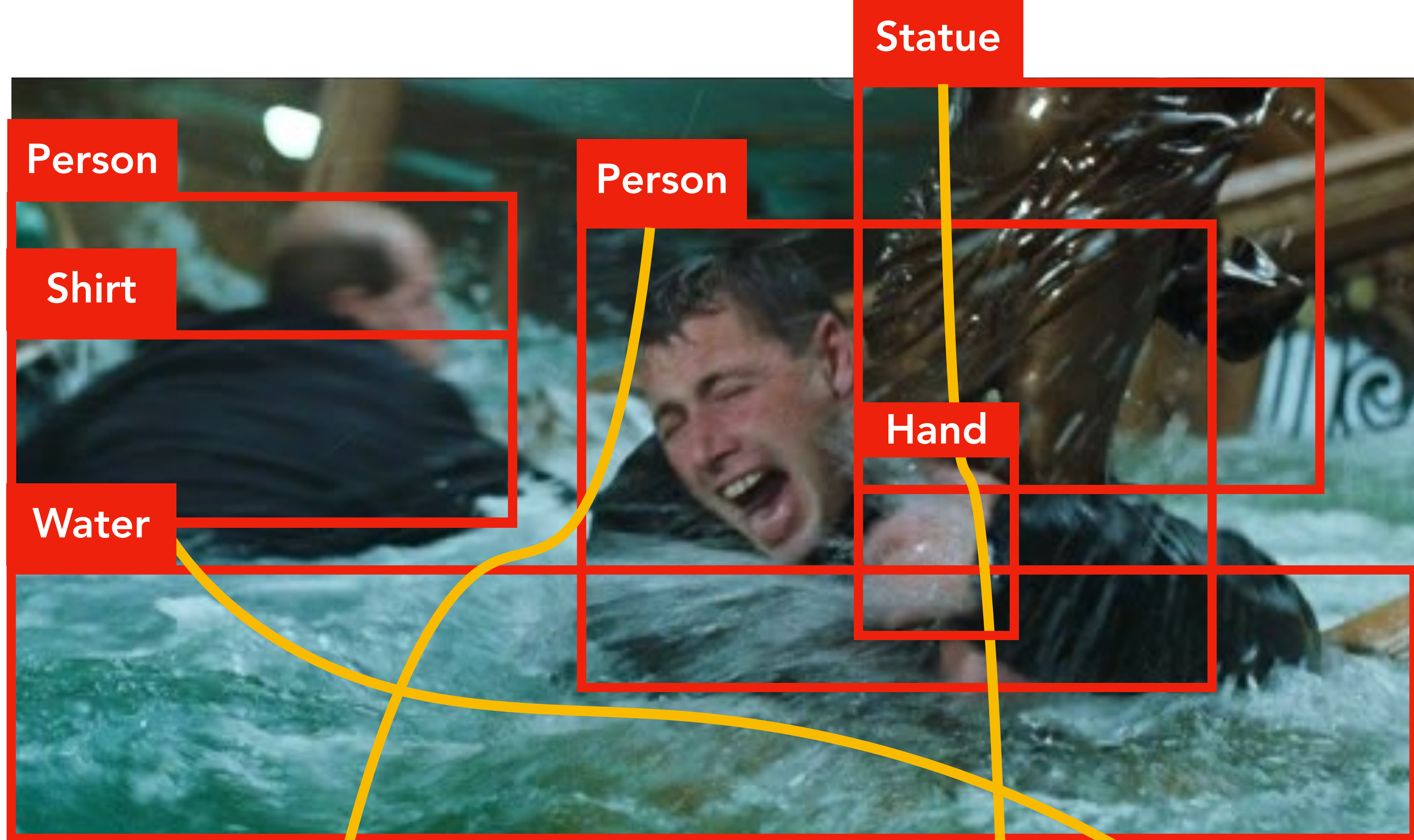




# Scene Graph

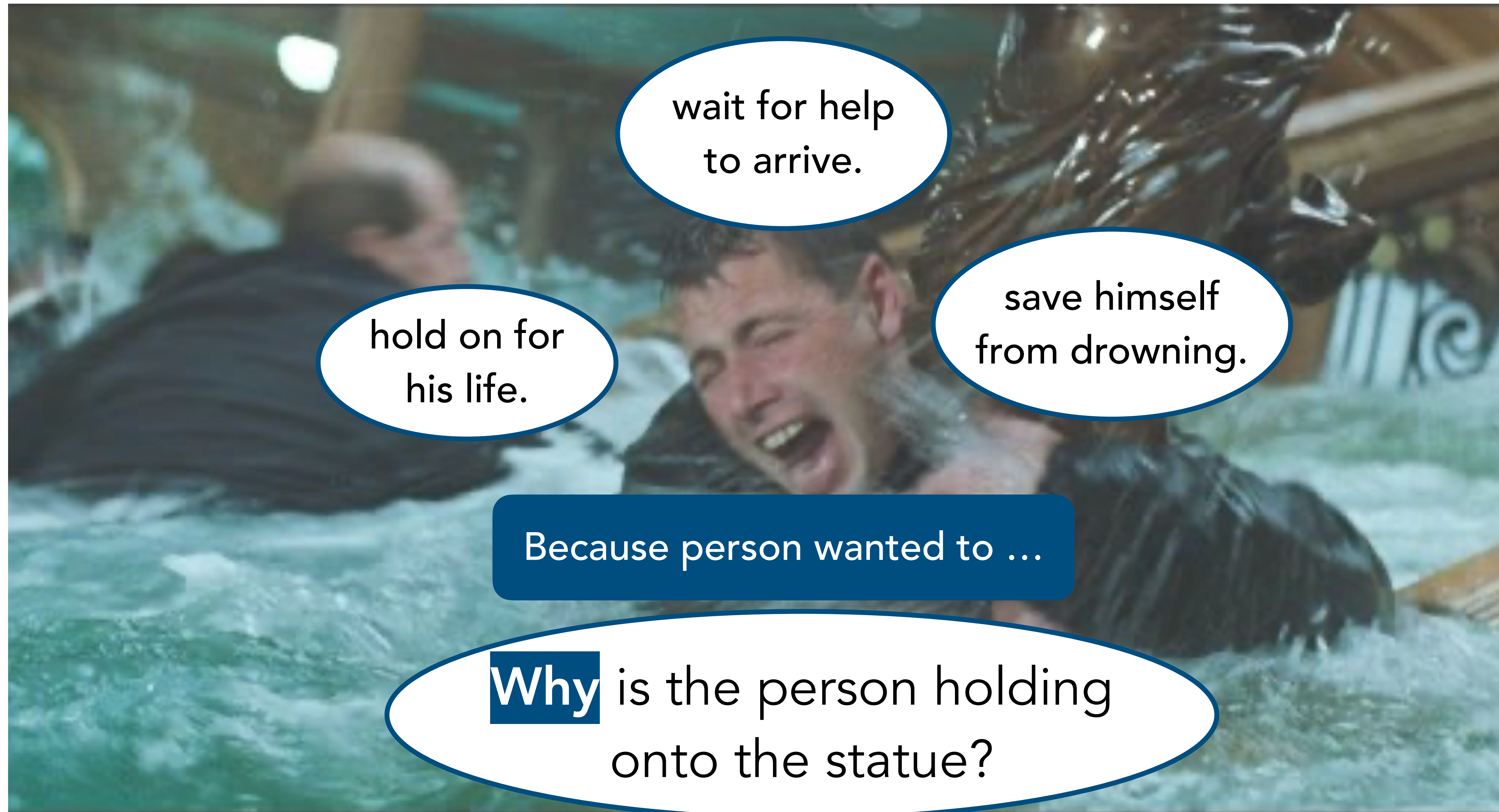
(Johnson et. al., 2015; Krishna et. al., 2016)





A person is holding onto a bronze statue in water.





wait for help  
to arrive.

hold on for  
his life.

save himself  
from drowning.

Because person wanted to ...

**Why** is the person holding  
onto the statue?

A person is holding onto a bronze statue in water.



What did person need to do **before** the image?

What will person do **after** the image?

Before, person needed to ...

Swim towards the statue.

Notice water washing in.

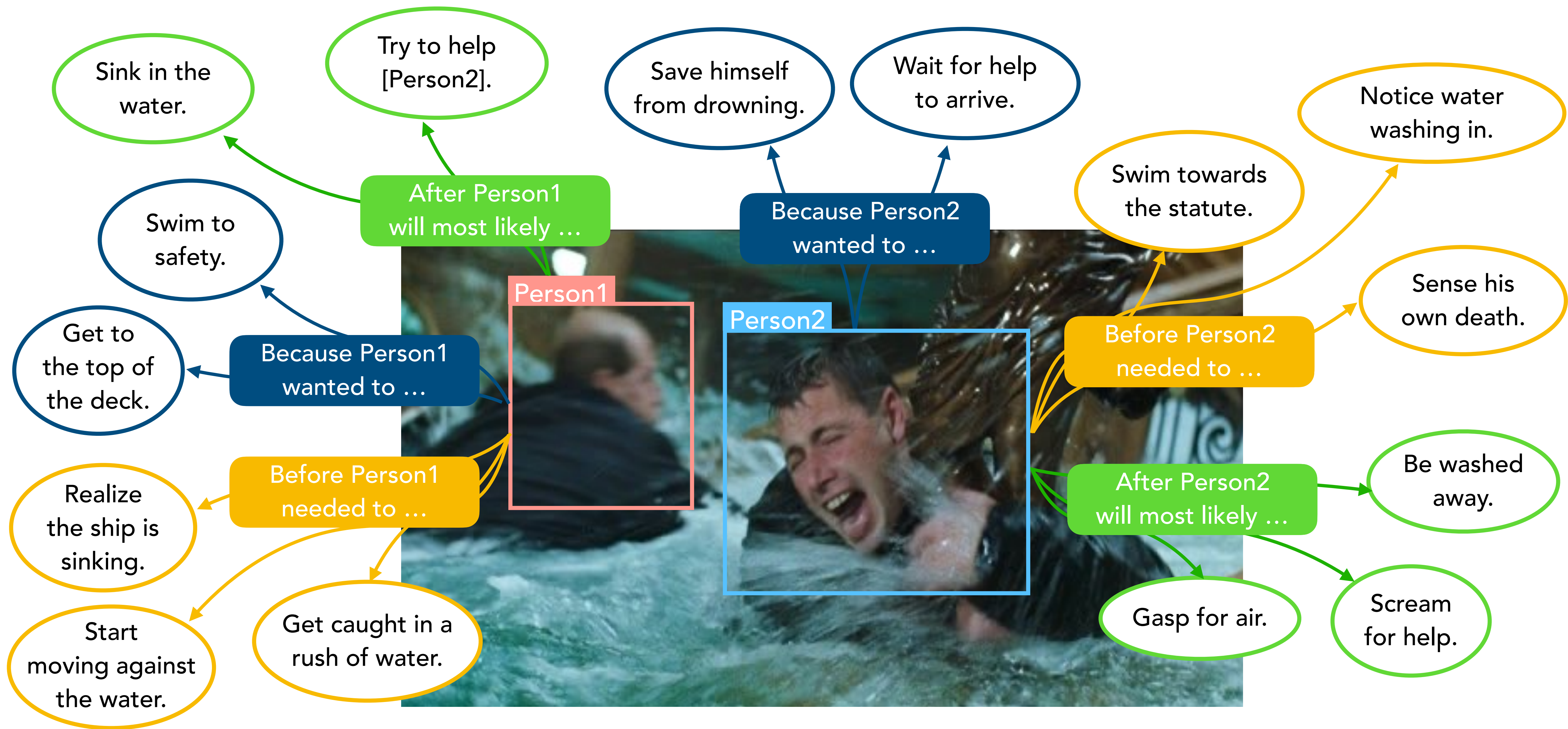
After, person will most likely ...

Gasp for air.

Be washed away.

A person is holding onto a bronze statue in water.





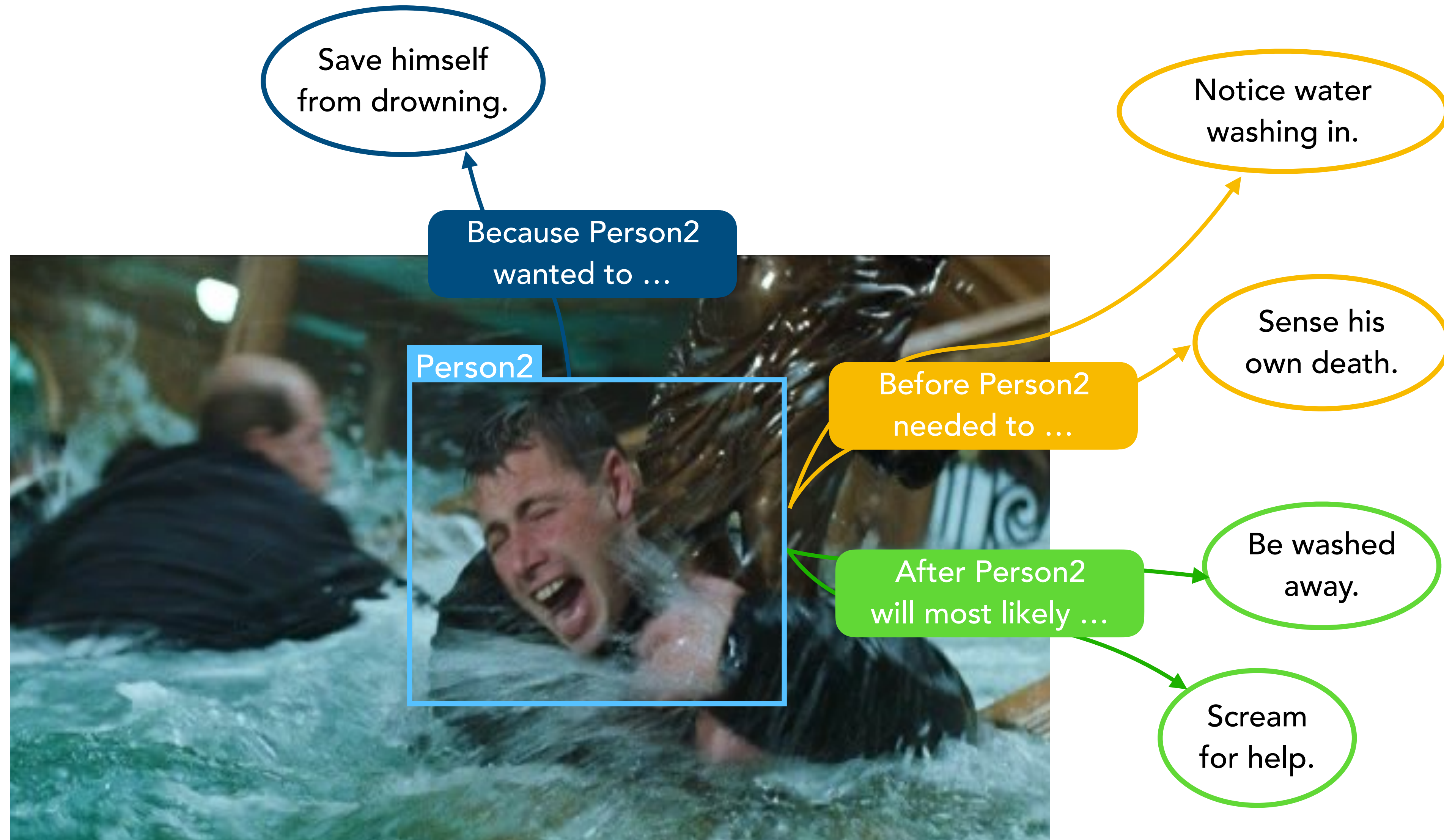






# Visual Commonsense Graphs:

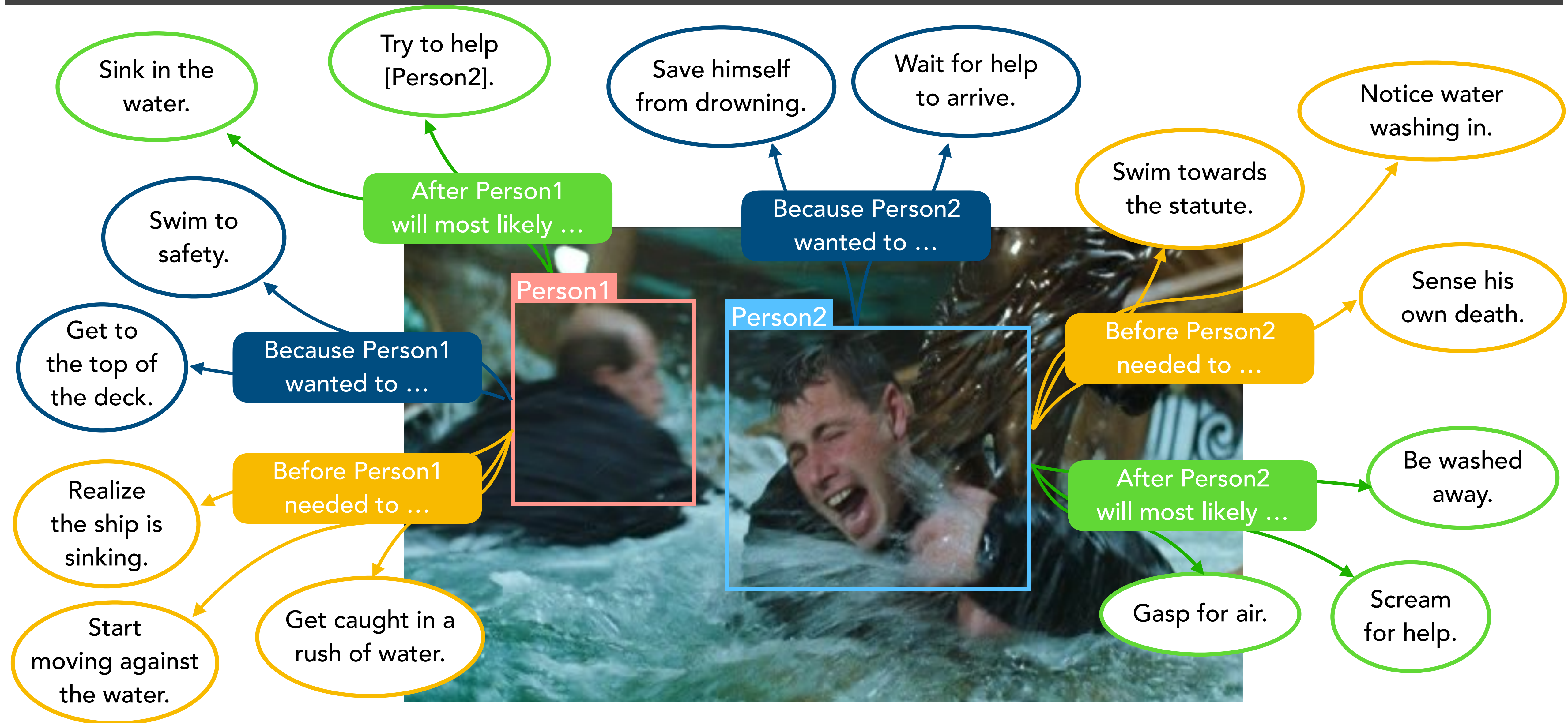
Reasoning about the *Dynamic* Context of a *Still* Image

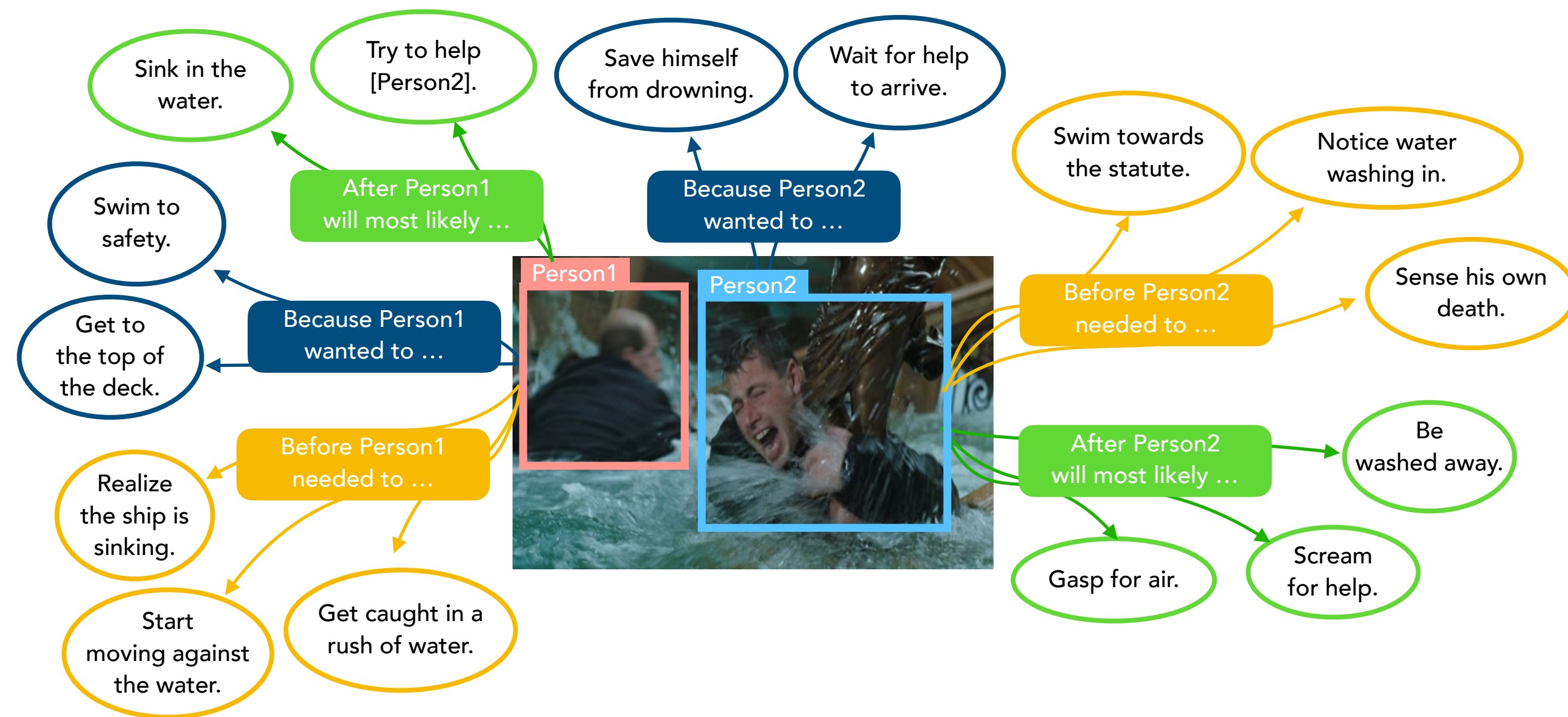




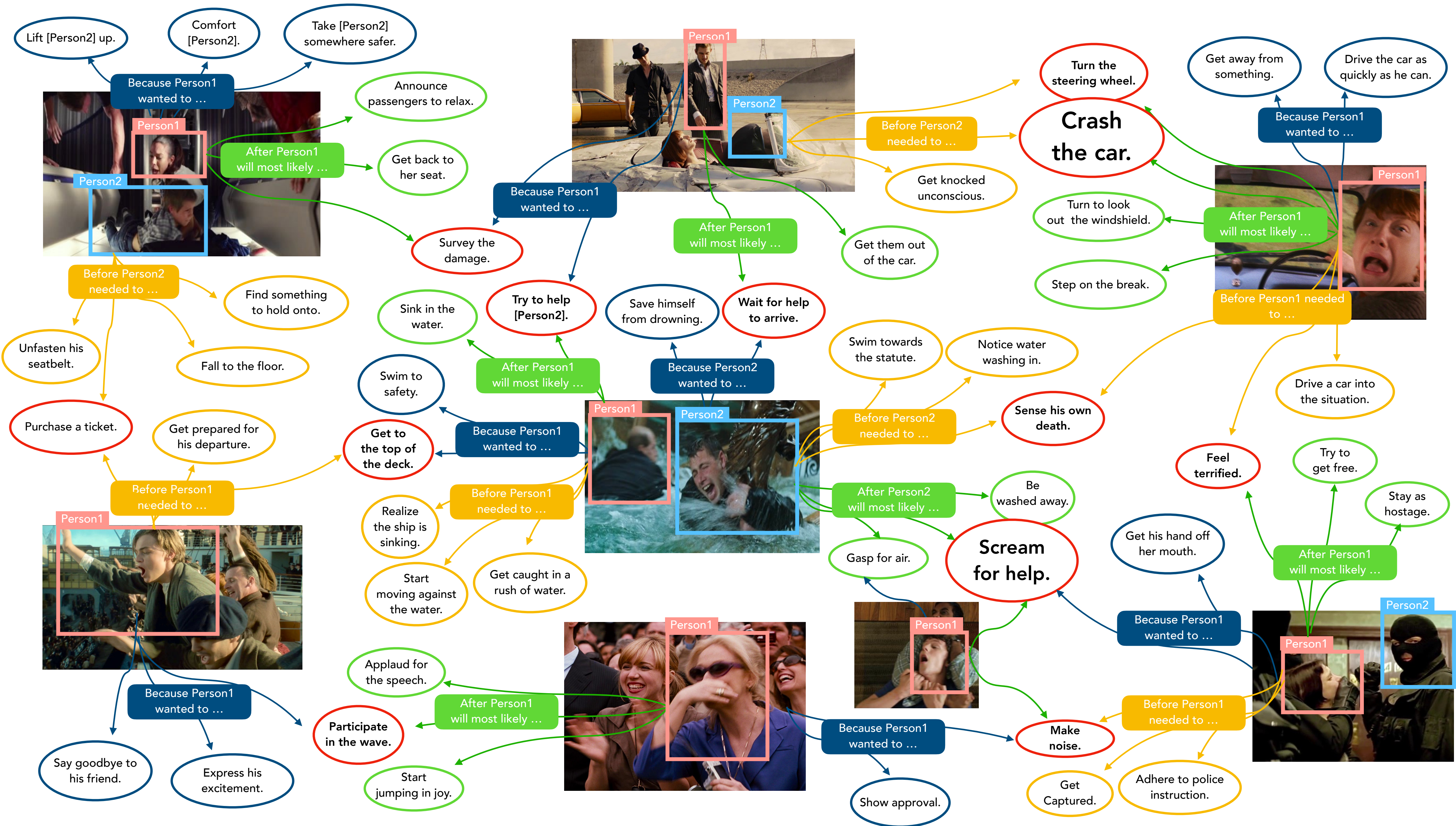
# Visual Commonsense Graphs:

Reasoning about the *Dynamic* Context of a *Still* Image

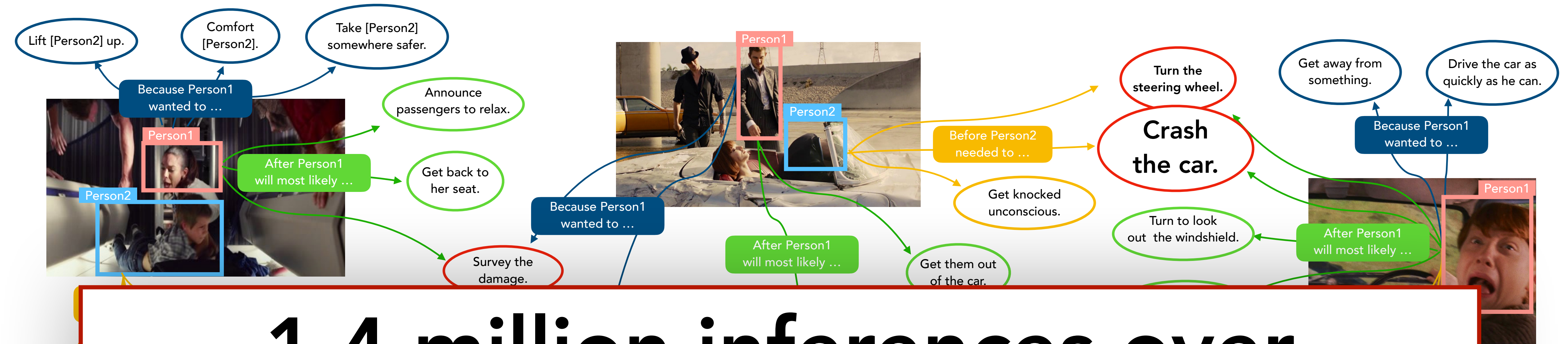












1.4 million inferences over

60K images with

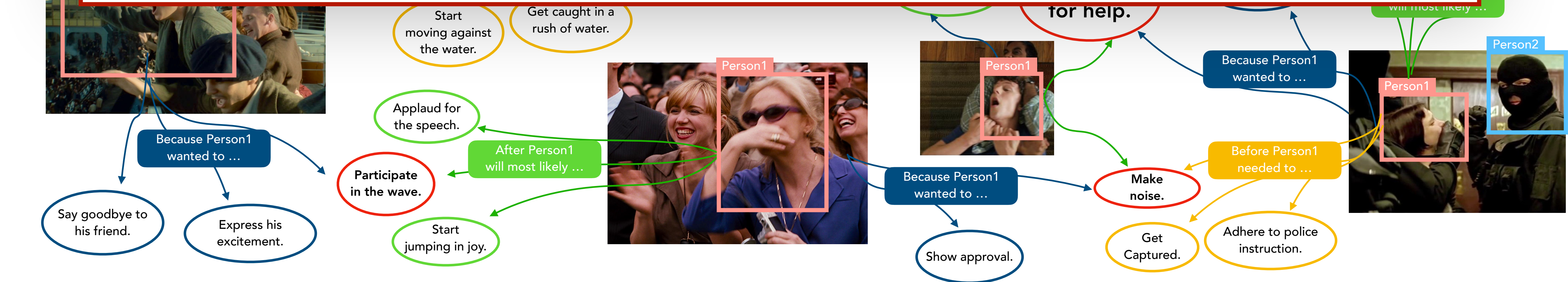
<https://visualcomet.xyz>

Unfasten h  
seatbelt.

Purchase

Pe

Stay as  
hostage.



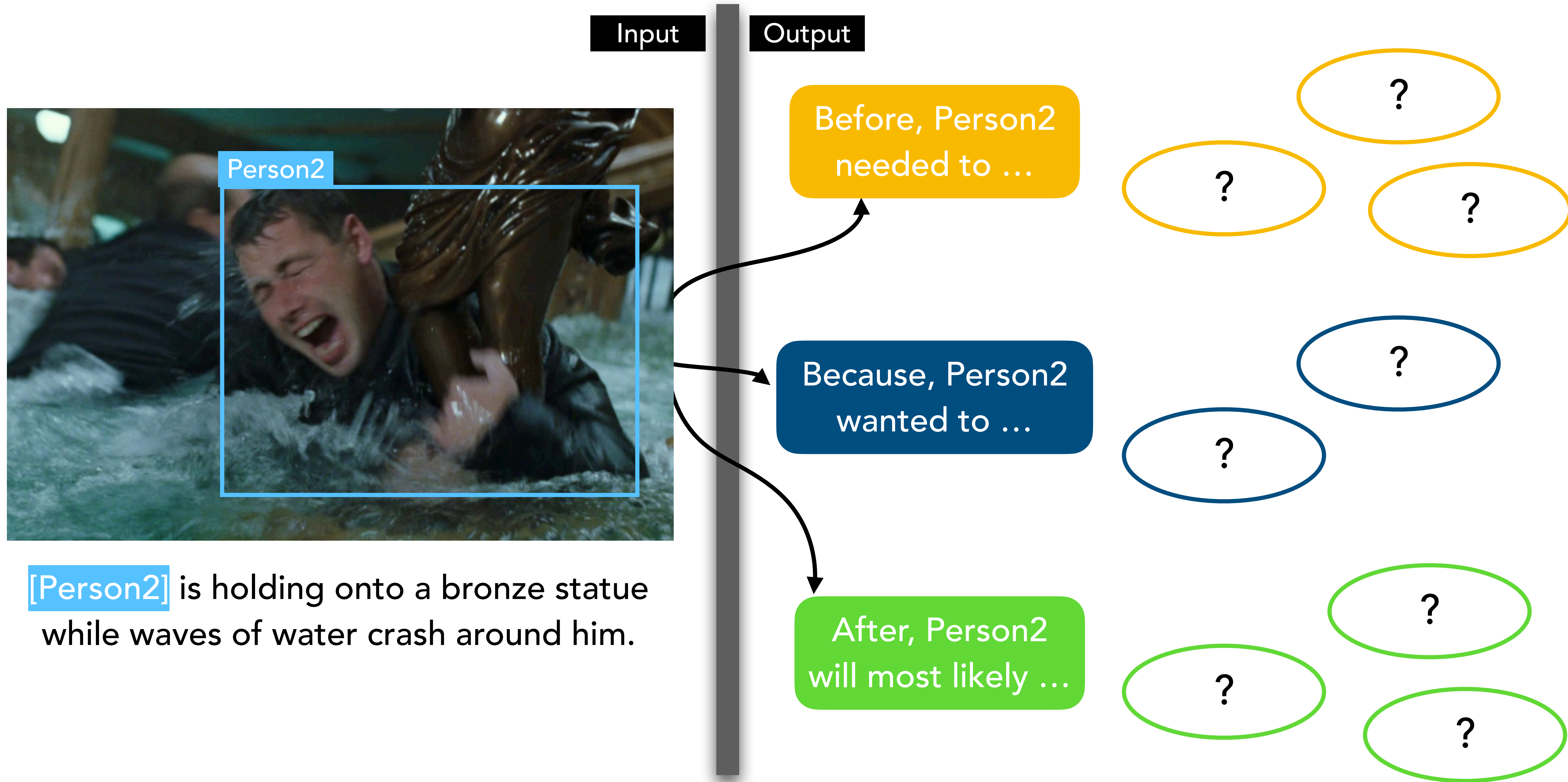


# Dataset Statistics

- **60K** Images from VCR (Zellers et al. 2019) complex scenes from movie scenes
- **130K** Person-Grounded **Event** Sentences (Min. 2 per Image)
- **1.4M** Person-Grounded Inference Sentences
  - Min. 4 **Before** Inferences per **Event**
  - Min. 4 **After** Inferences per **Event**
  - Min. 2 **Intent** Inferences per **Event**
- All Sentences Annotated with Amazon Mechanical Turk

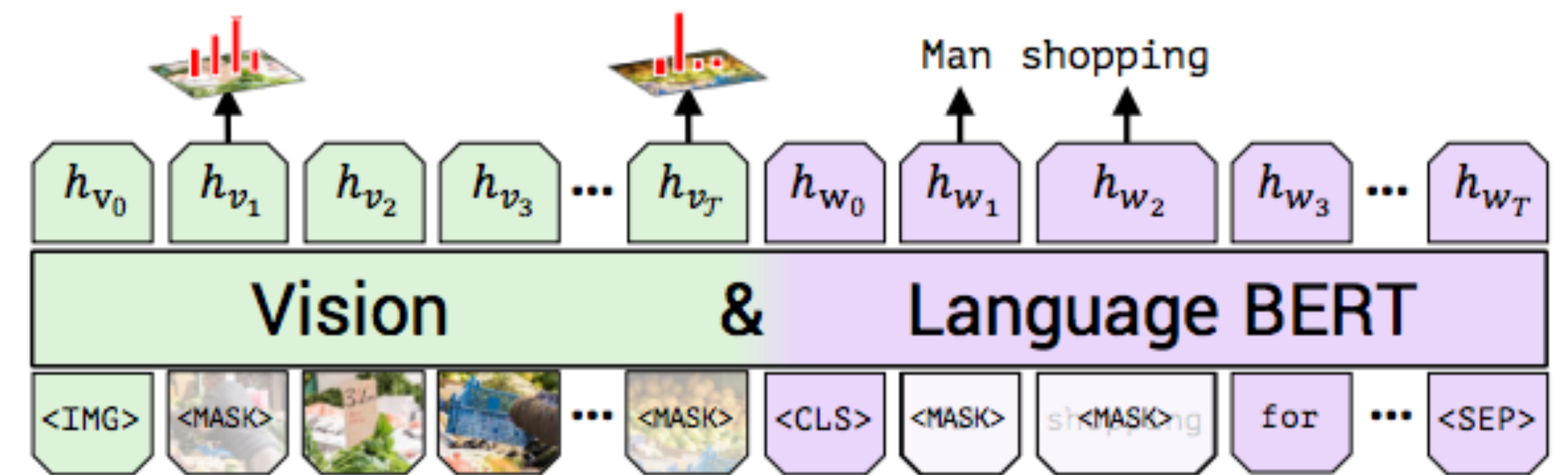
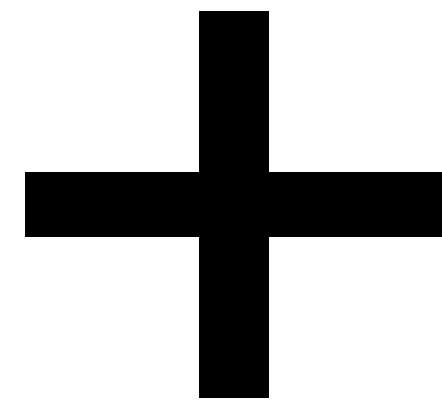
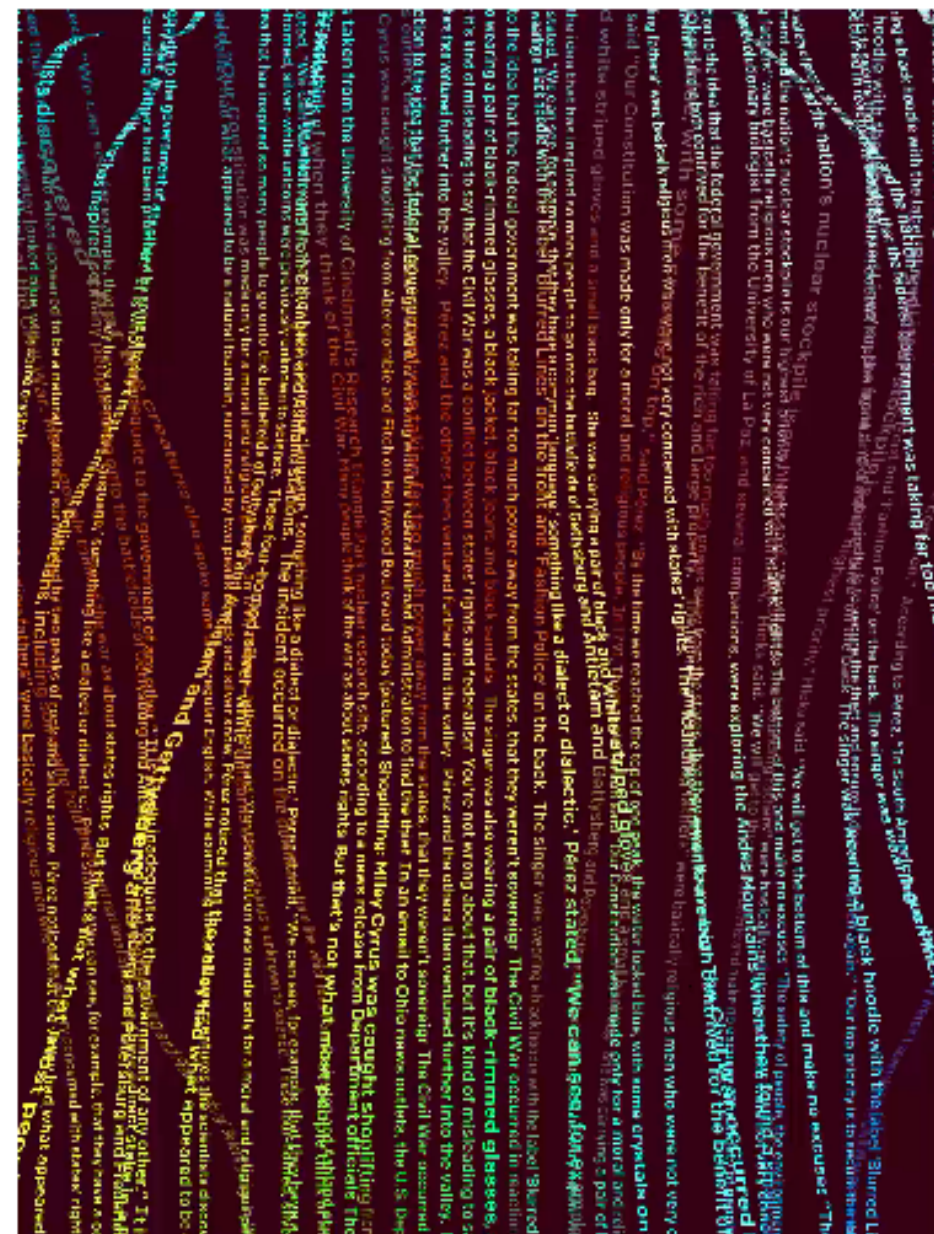
Annotators had access  
to before/after videos.

# Task: Generating Commonsense Inferences in Language





# Our Model Builds on Pre-Trained Language Models



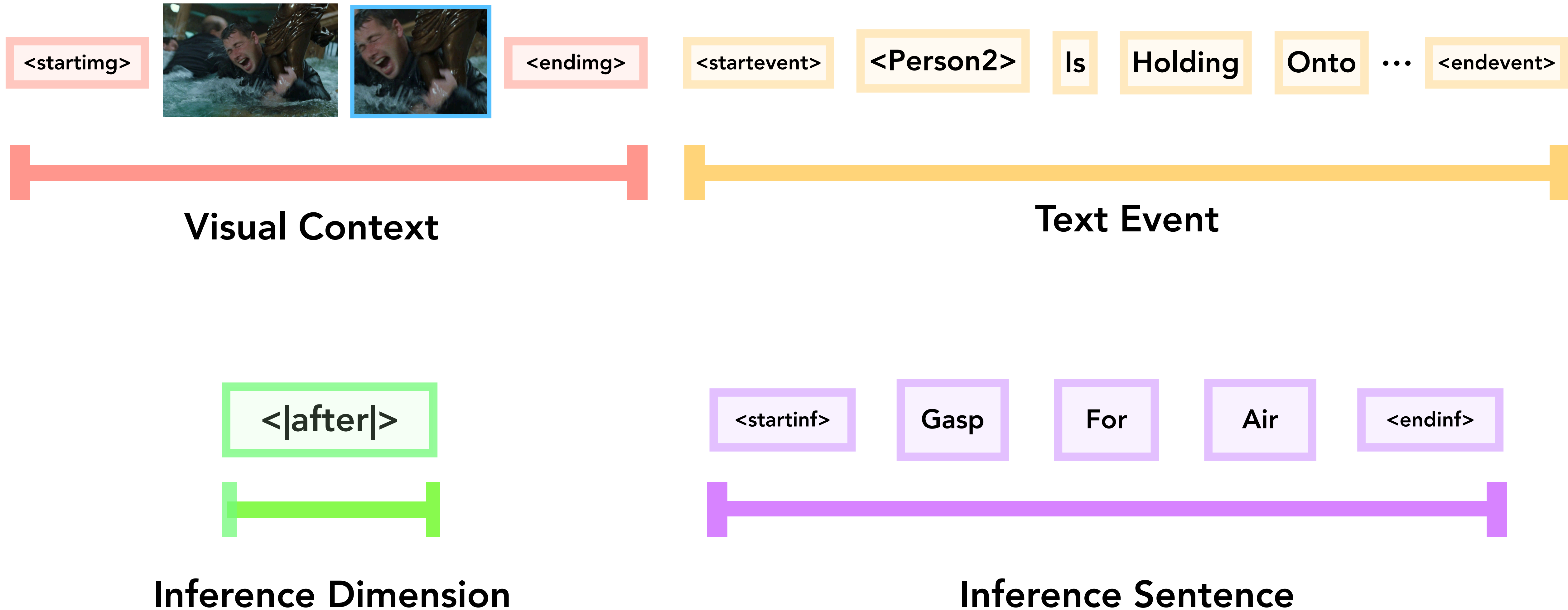
**GPT-2 for Conditional  
Generation**

*(Radford et. al., 2019)*

**Vision-Language  
Transformer Architecture**

*(Lu et. al., 2020; Su et. al, 2020; Tan et. al, 2020)*

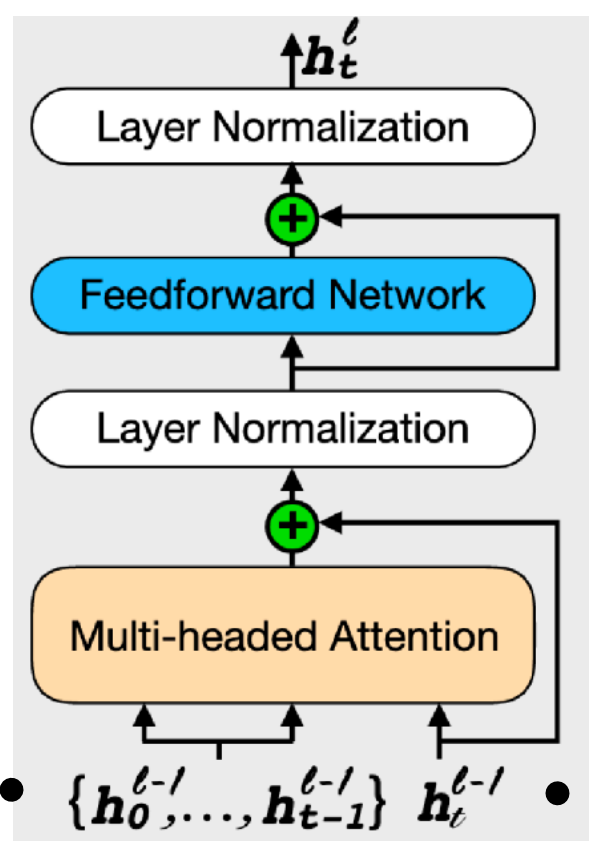
# Our Approach



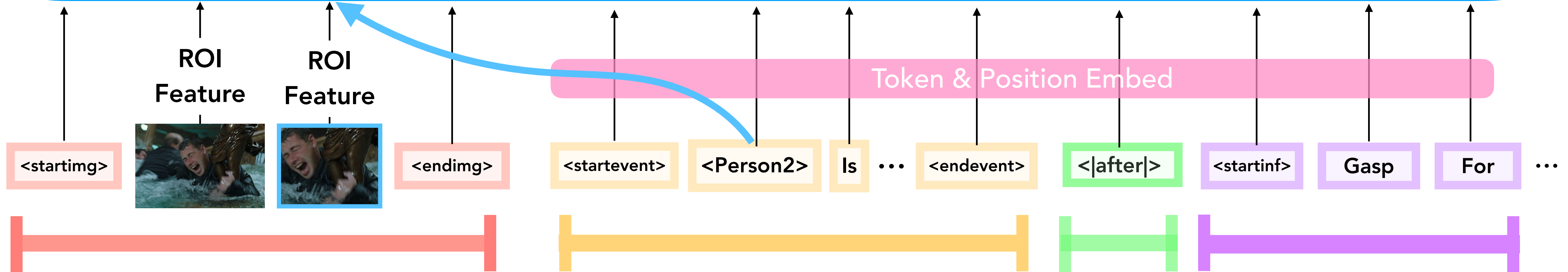
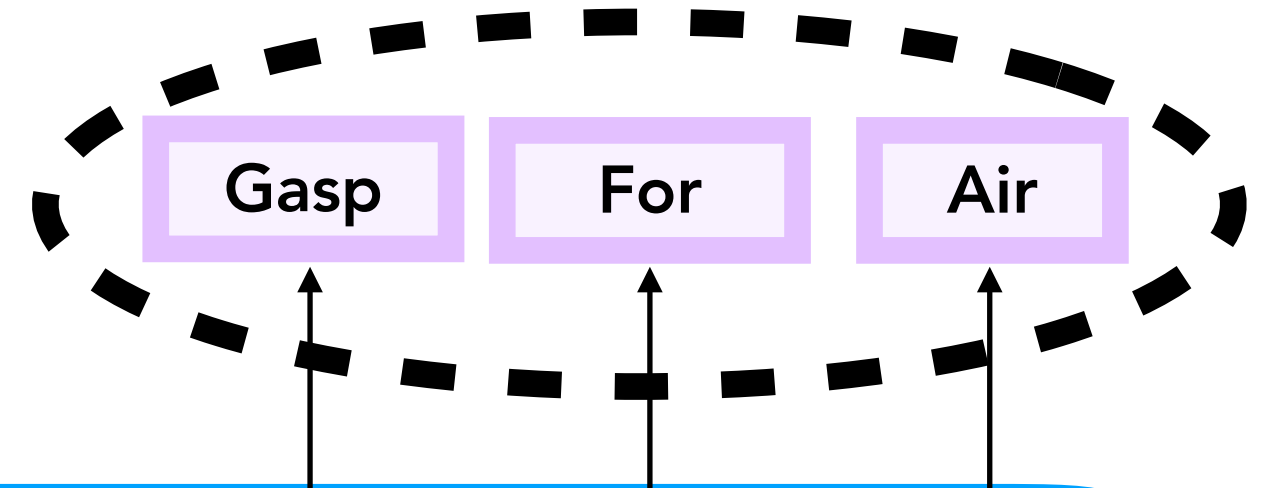


# Model Architecture

(Vaswani et al., 2017)



Fine-tune LM:



Visual Context

Text Event

Inference Dim. & Sentence

# Sanity Checks against Spurious Dataset Biases

In **MS CoCo**, captions retrieved from KNN matched human-level performance

*Devlin et al. 2015*

=> **Sanity Check #1:**  
Nearest Neighbor Baseline

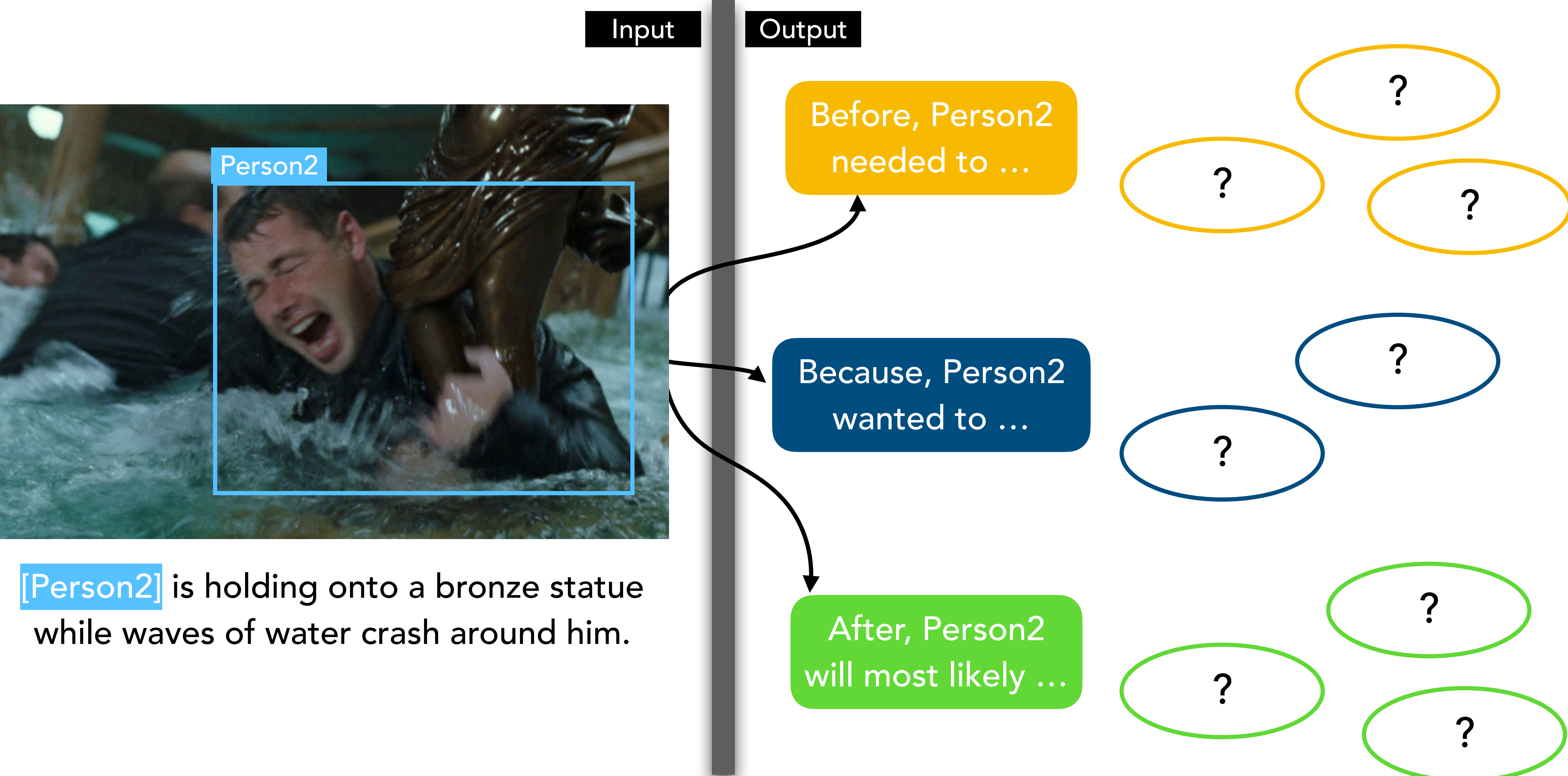
In (the original) **VQA**, models did too well without looking at the image:

1. "what is the color of the banana?"
  - "yellow"
2. "how many ...?"
  - "two"



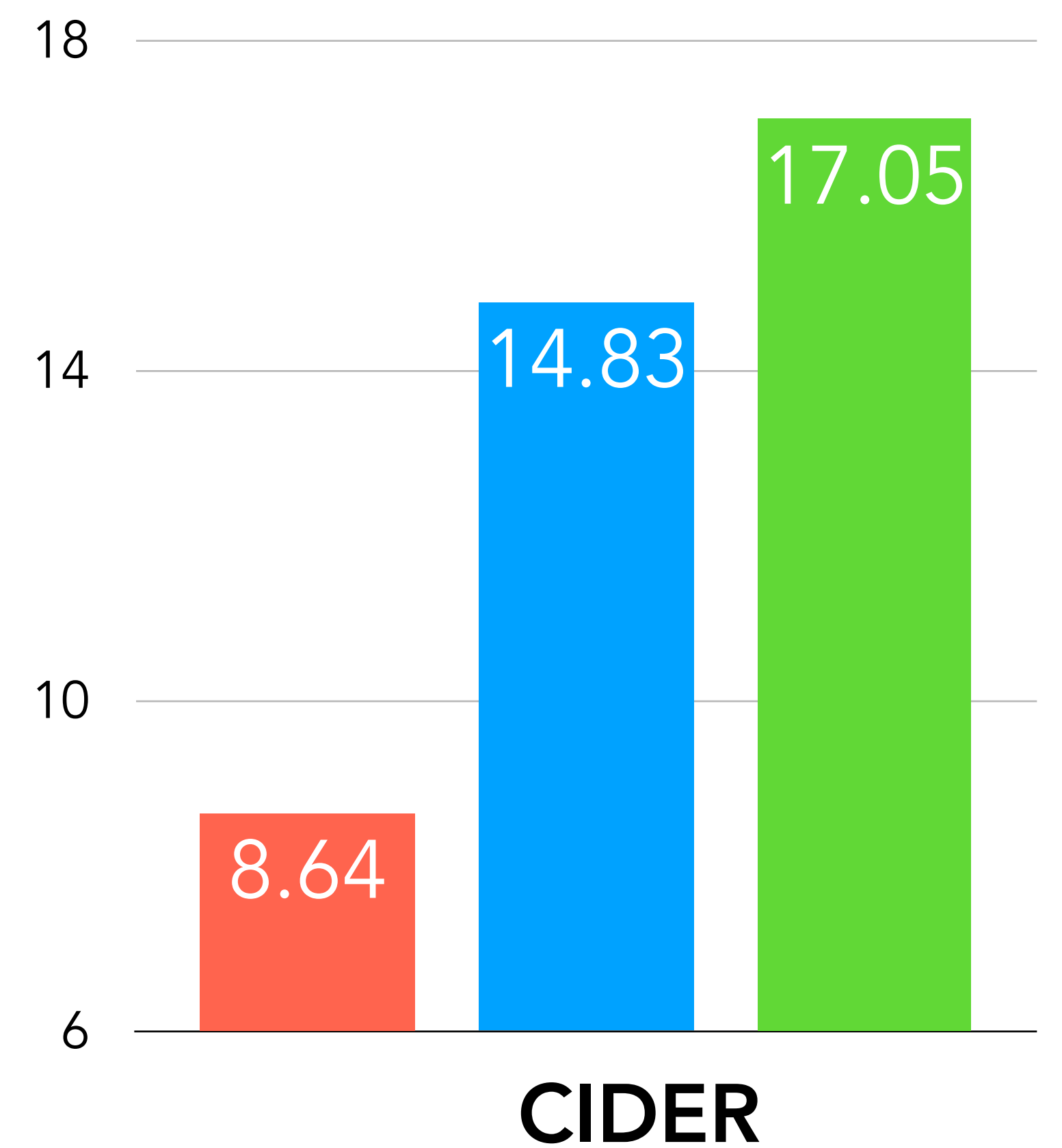
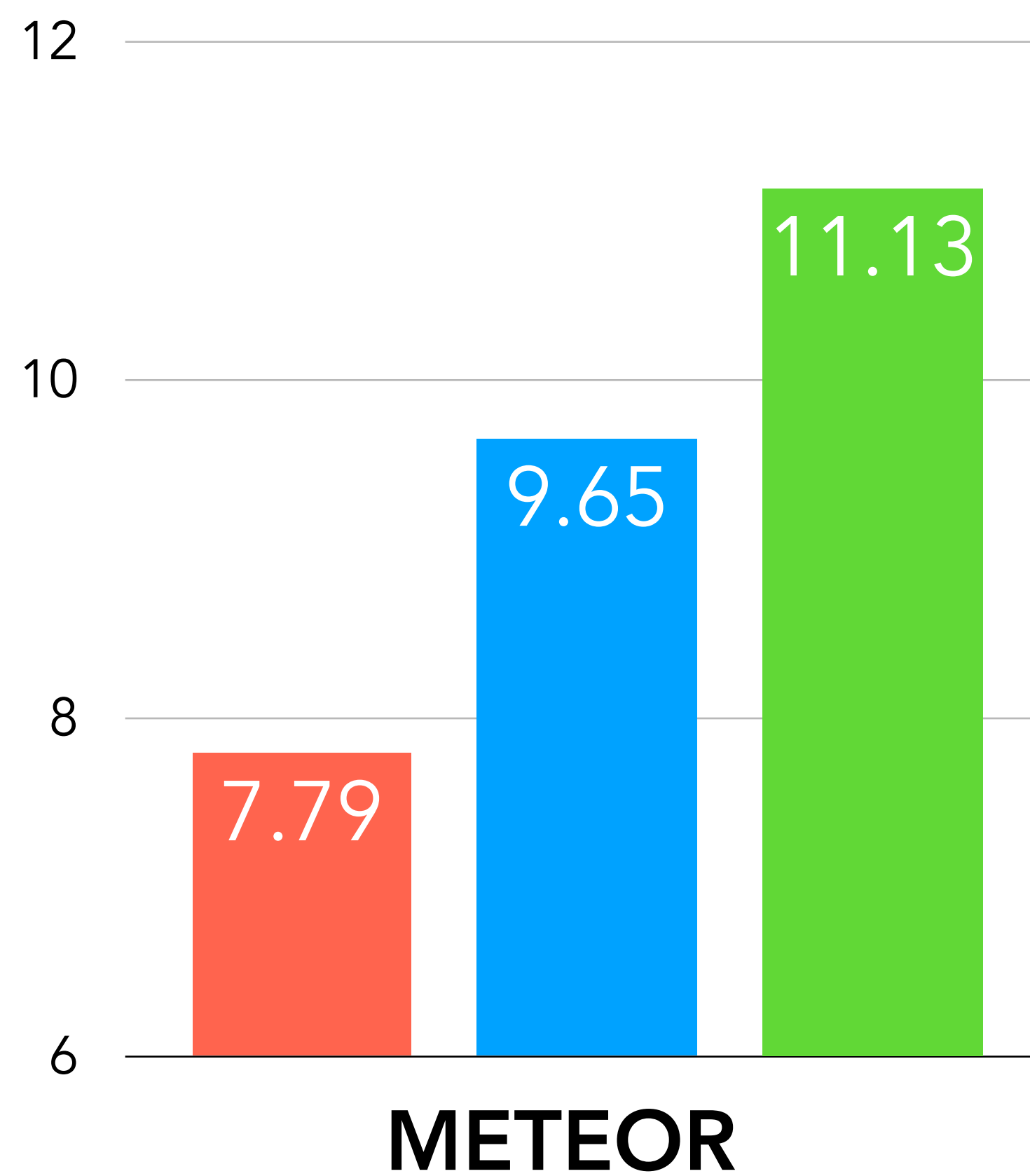
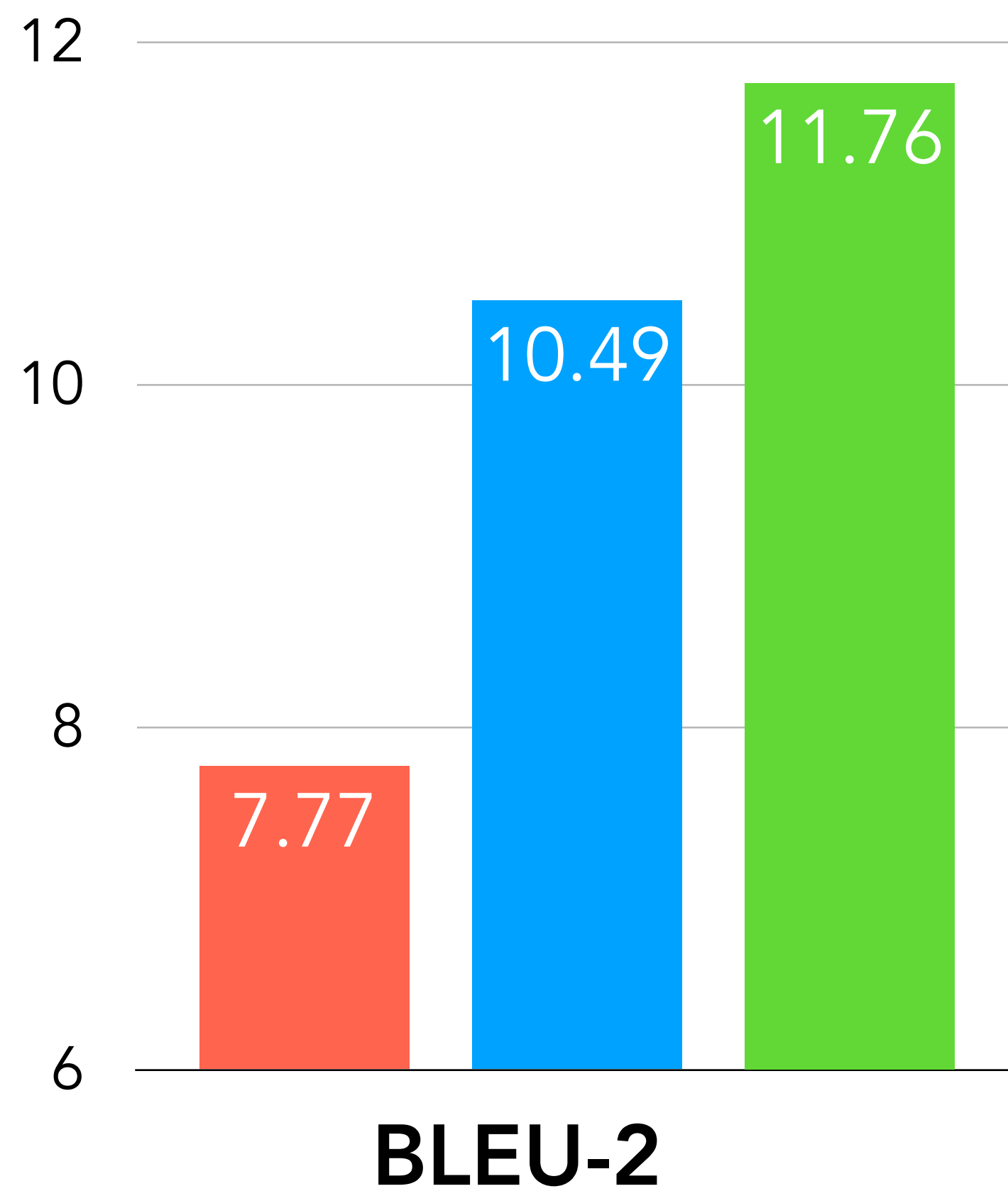


# Sanity Check #2: Language-Only Baseline Ignoring the Image



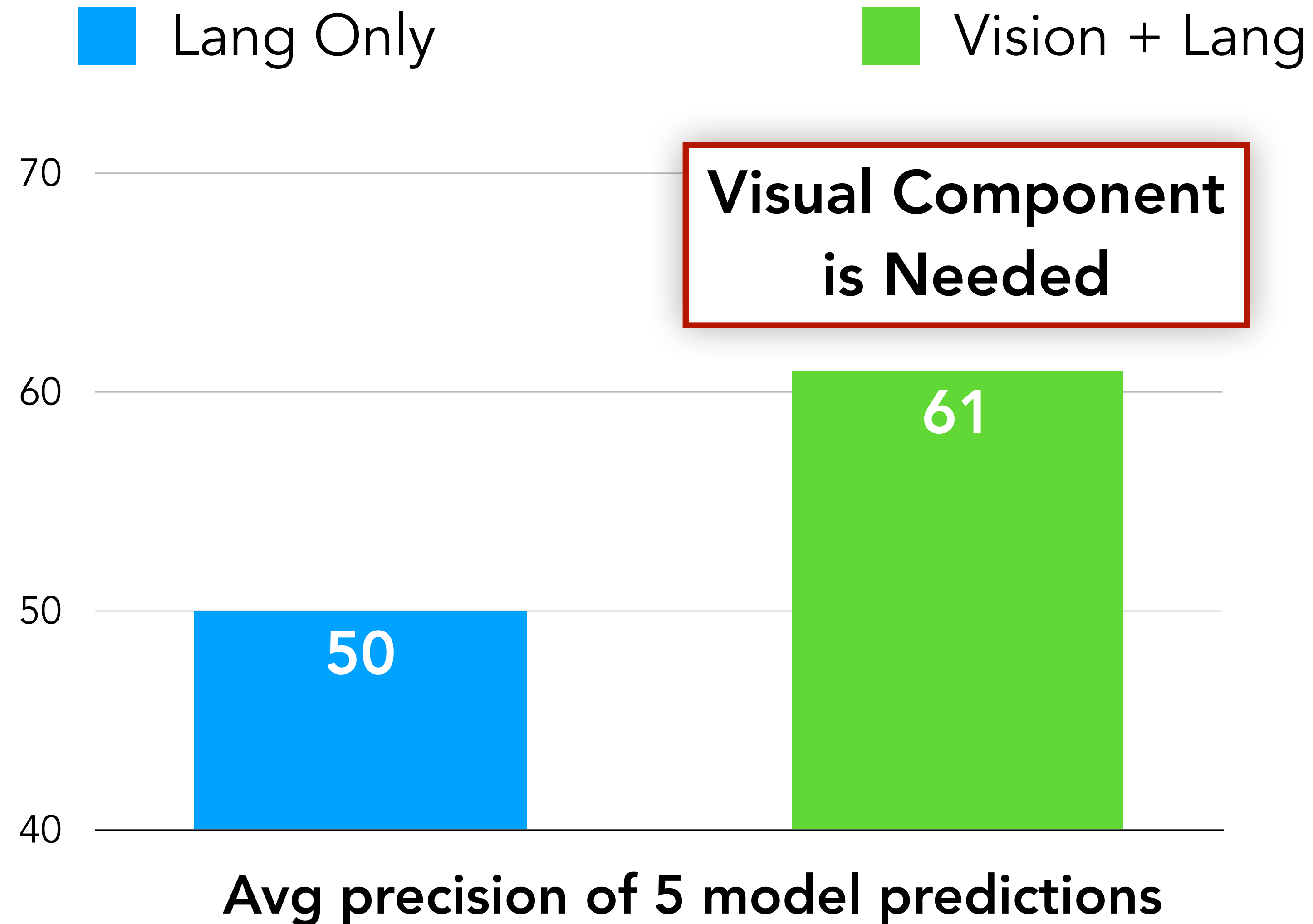
# Automatic Evaluation

■ Nearest Neighbor ■ Lang Only ■ Vision + Lang





# Human Evaluation



Before, Person1 needed to ...

Unlikely

Input

Output

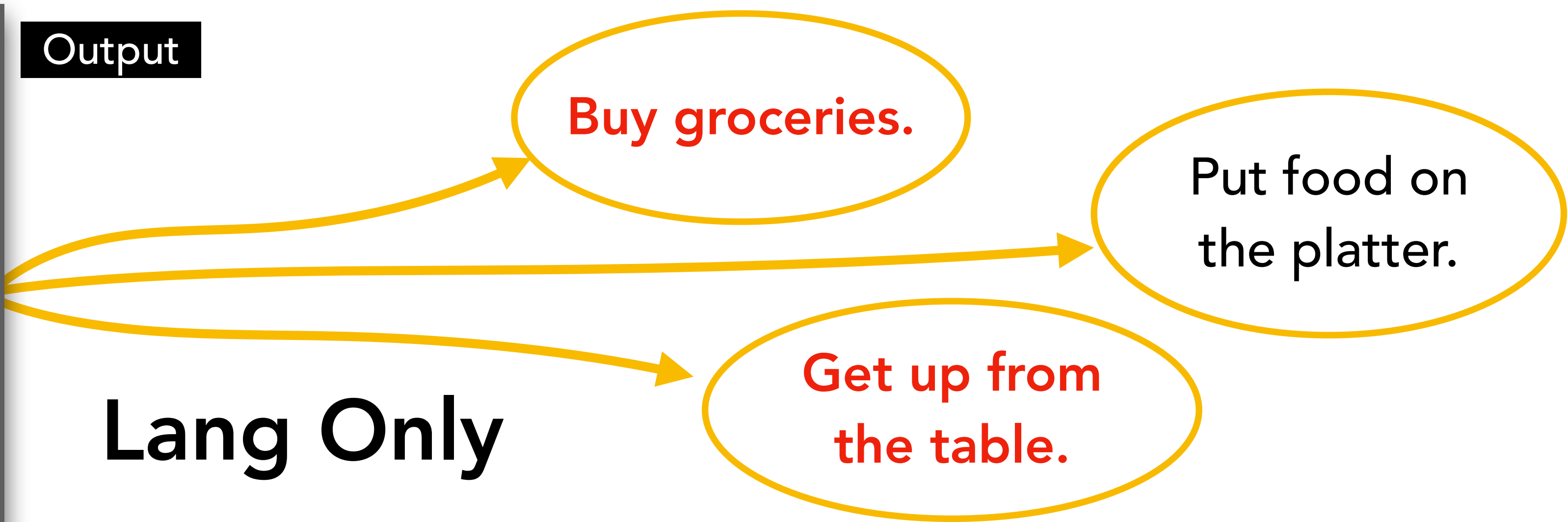
[Person1] is putting a platter on the table at an outdoor restaurant.

Buy groceries.

Put food on the platter.

Get up from the table.

Lang Only





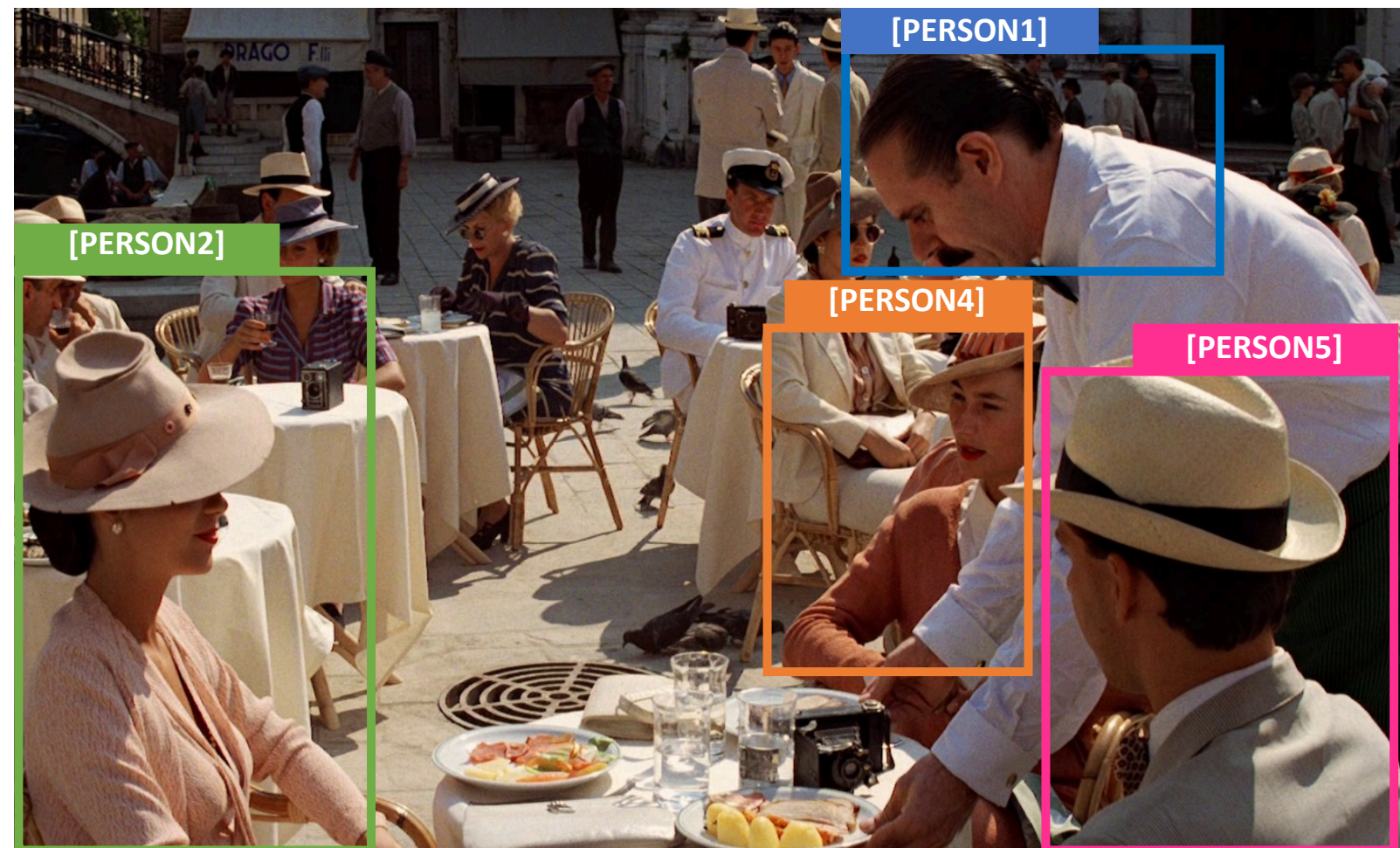
Before, Person1 needed to ...

Unlikely

Input

Output

[Person1] is putting a platter on the table at an outdoor restaurant.



Lang Only

Buy groceries.

Put food on the platter.

Get up from the table.

Vision + Lang

Be hired as a waiter.

Wait for everyone to sit down.

Receive an order for platter.



Because, Person1 wanted to ...

Unlikely

Input

Output

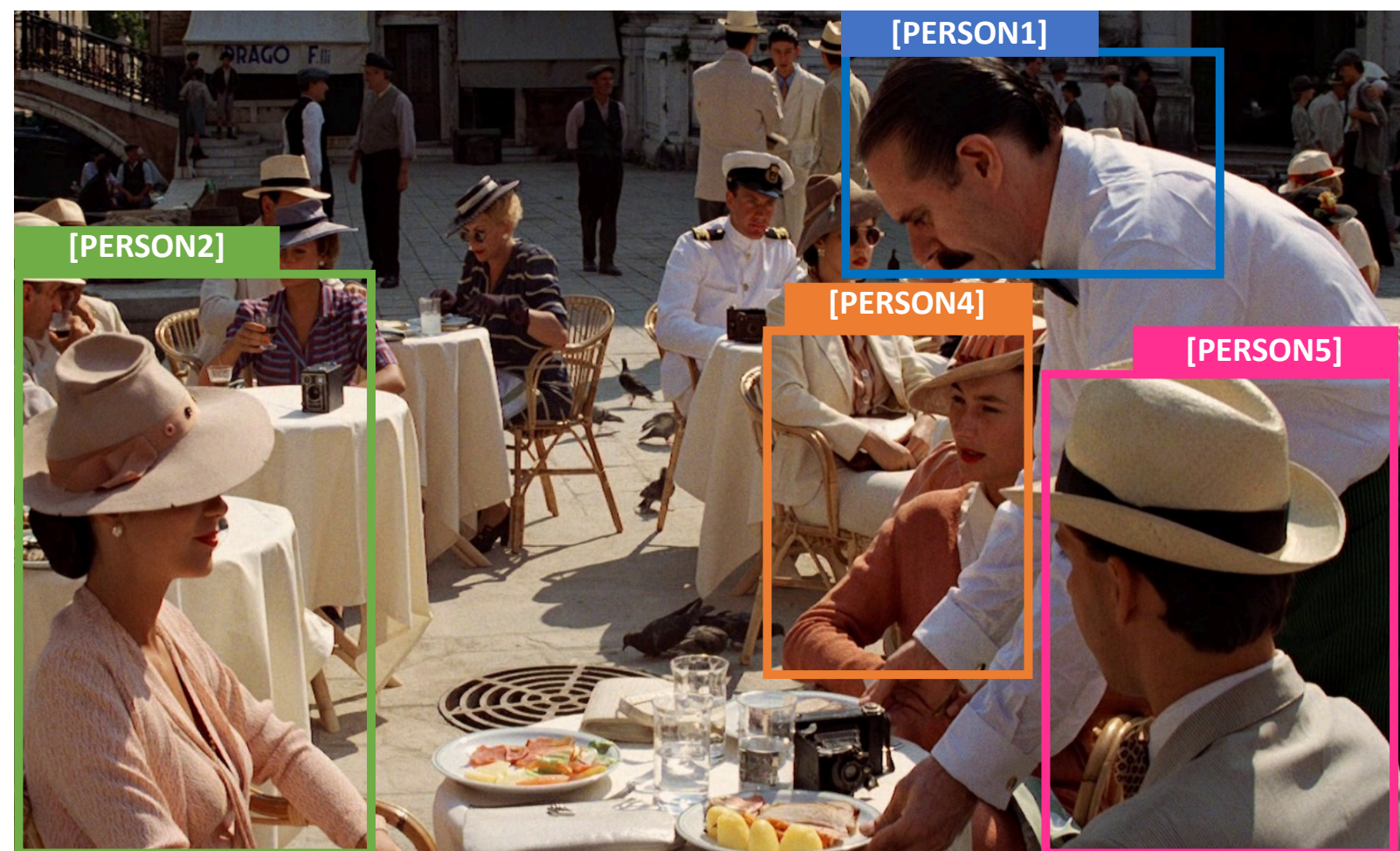
[Person1] is putting a platter on the table at an outdoor restaurant.

Have dessert

Tend to the patrons.

Ensure the food is taken care of.

Lang Only





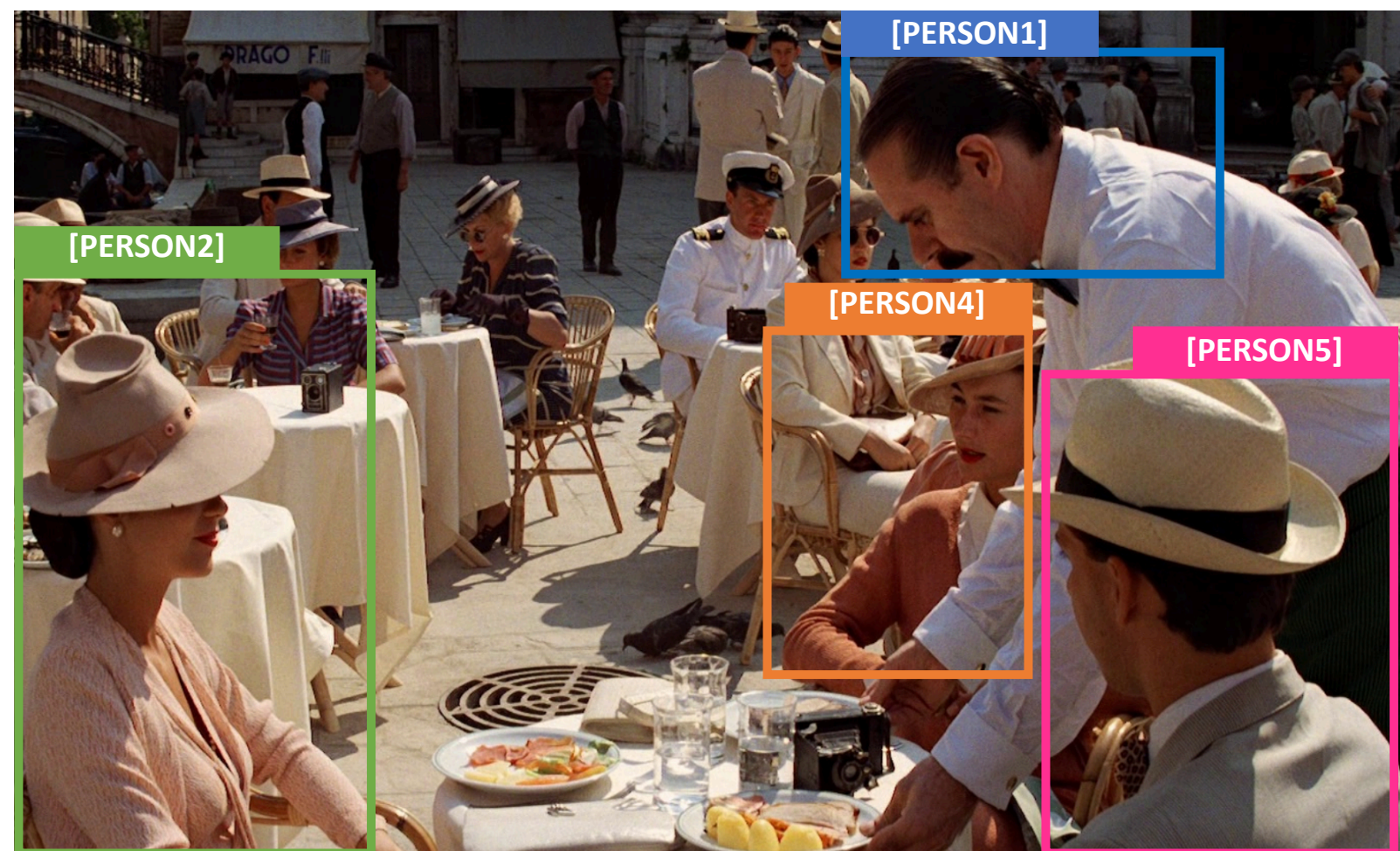
Because, Person1 wanted to ...

Unlikely

Input

Output

[Person1] is putting a platter on the table at an outdoor restaurant.



Lang Only

Have dessert

Tend to the patrons.

Ensure the food is taken care of.

Vision + Lang

Serve [P2], [P4], and [P5].

Greet [P2], [P4], and [P5].

Have [P2], [P4], and [P5] to eat.



After, Person1 will most likely ...

Unlikely

Input

Output

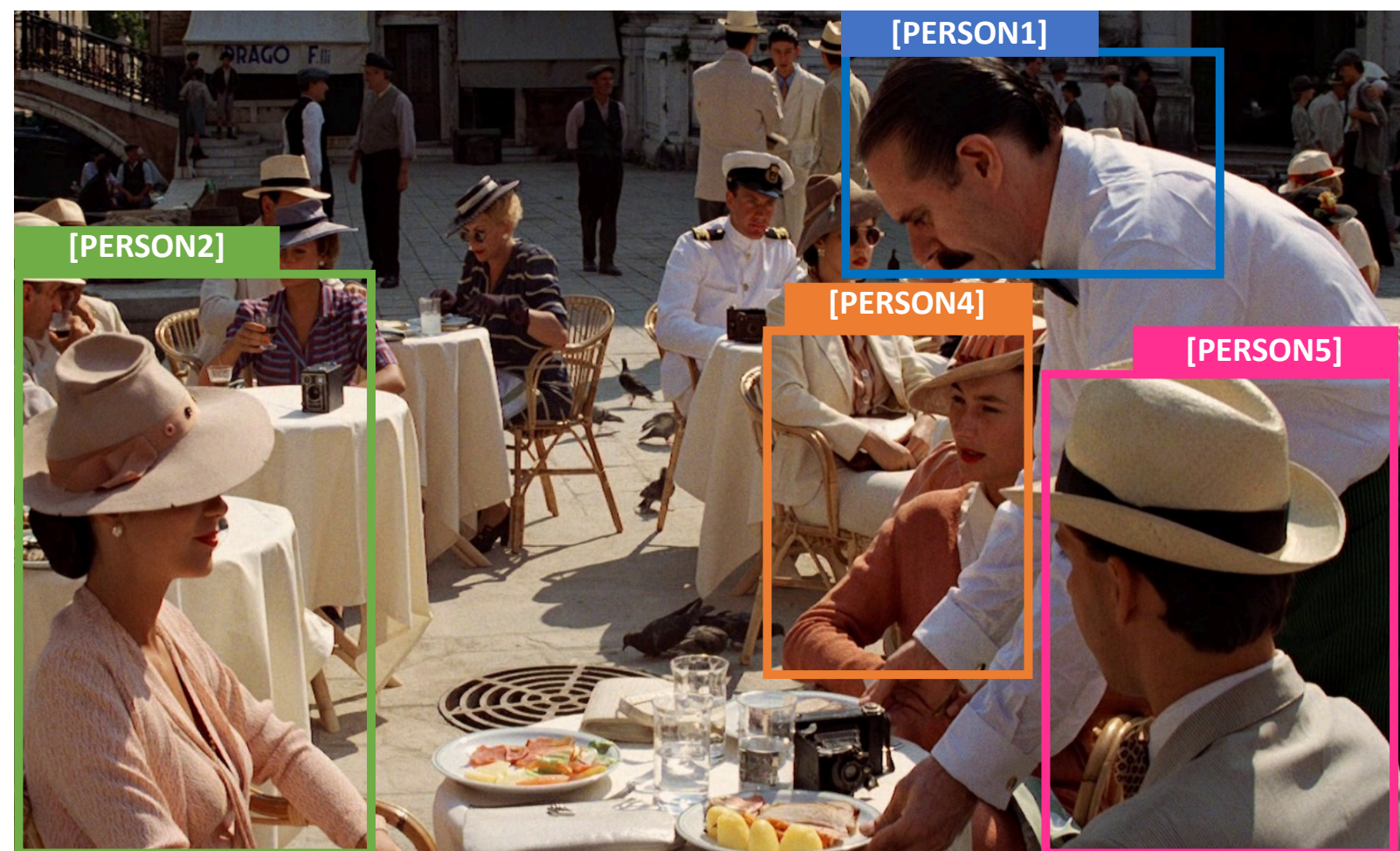
[Person1] is putting a platter on the table at an outdoor restaurant.

Sip the water.

Ask [P2] for a menu.

Get up and walk over to his table.

Lang Only





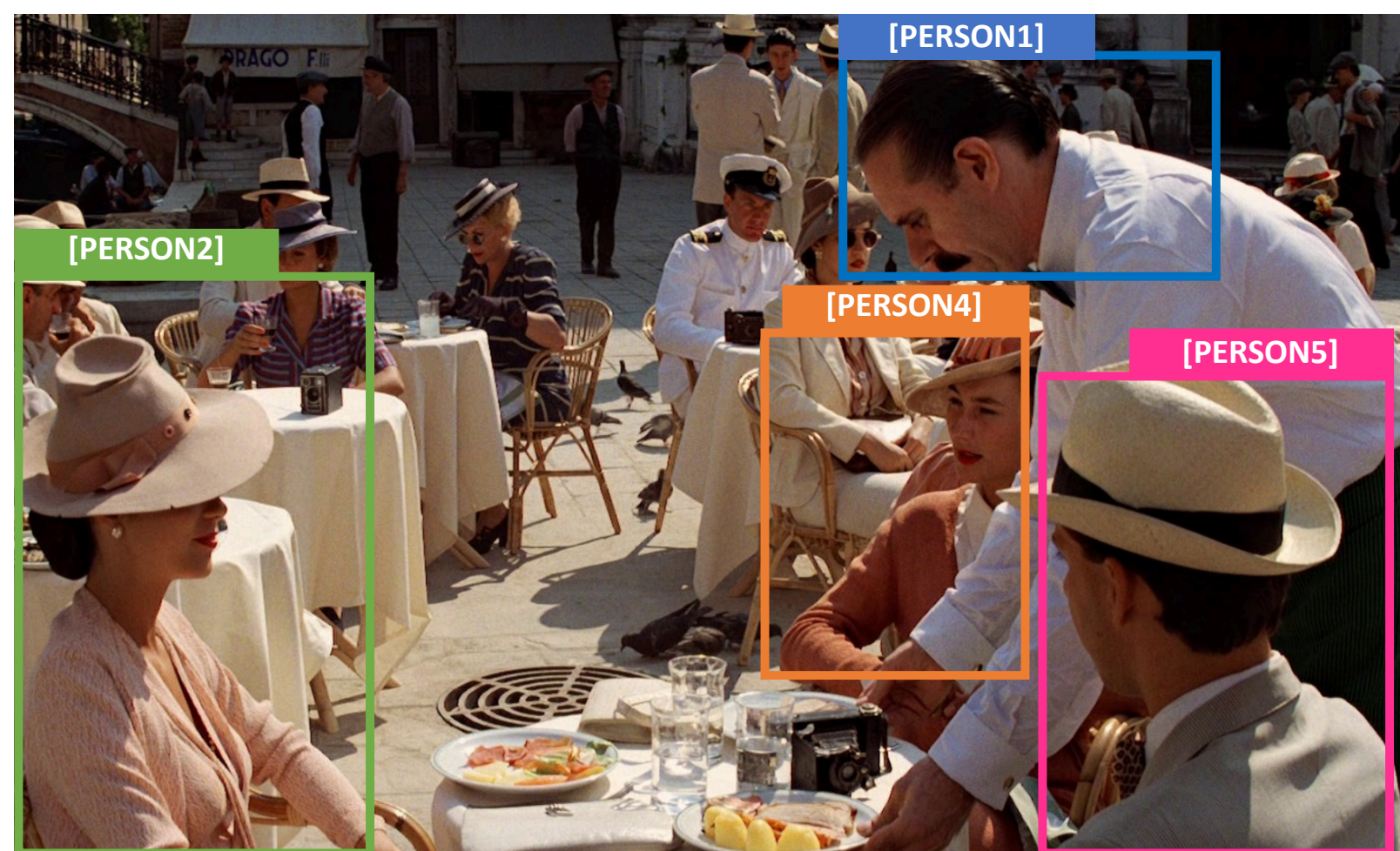
After, Person1 will most likely ...

Unlikely

Input

Output

[Person1] is putting a platter on the table at an outdoor restaurant.



Lang Only

Sip the water.

Ask [P2] for a menu.

Get up and walk over to his table.

Vision + Lang

Take drinks.

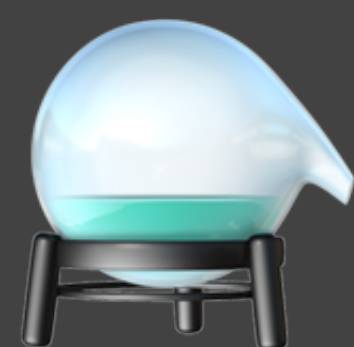
Get back to his work duties.

Get back to the kitchen to get more food.

# Conclusion

- **Visual Commonsense Graphs**, the first large-scale repository of visual commonsense with more than **1.4 million commonsense inferences** over **60k images with complex visual scenes**.
- Promising results that use the power of pre-trained language to close the gap between **perception-level** and **cognition-level** visual understanding.
- The dataset supports new tasks and models beyond what are presented here.





# Social Chemistry 101

Learning to Reason about Social and Moral Norms

EMNLP 2020

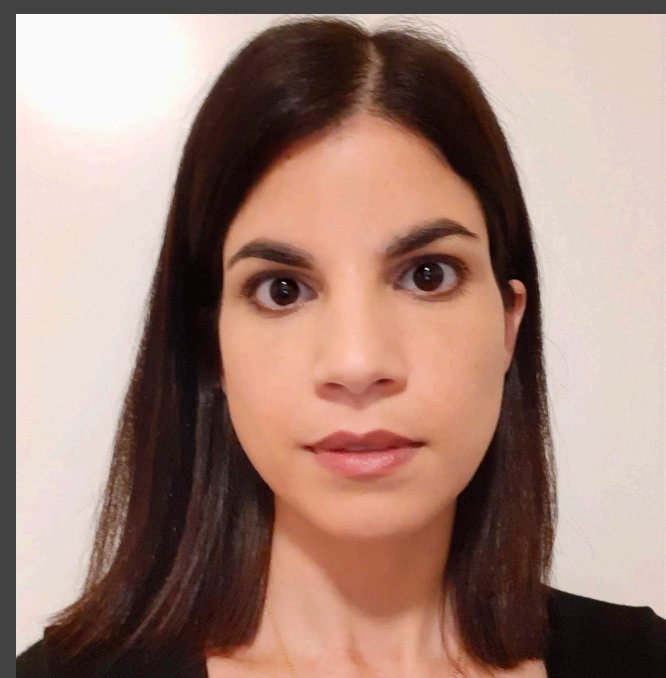
Maxwell Forbes



Jena Hwang



Vered Shwartz



Maarten Sap



Me





# GPT3's Morality?

- Running a blender at 5am is rude because...
  - you can wake up the entire neighborhood. You can only do it if you are making a thick smoothie and need to incorporate some ice. 😊

Need to teach AI with people's descriptive ethics (declarative knowledge) as brute-force large networks (& data) won't cut it.

- It is ok to post fake news if...
  - it helps the xxx agenda, even if it hurts the country. 😱

Can't do social & moral norms without the full scope of language!

## SOCIAL CHEMISTRY 101: Learning to Reason about Social and Moral Norms

Maxwell Forbes<sup>†‡</sup> Jena D. Hwang<sup>‡</sup> Vered Shwartz<sup>†‡</sup> Maarten Sap<sup>†</sup> Yejin Choi<sup>†‡</sup>

<sup>†</sup>Paul G. Allen School of Computer Science & Engineering, University of Washington

<sup>‡</sup>Allen Institute for AI

{mbforbes, msap, yejin}@cs.washington.edu, {jenah, vereds}@allenai.org

[maxwellforbes.com/social-chemistry](https://maxwellforbes.com/social-chemistry)

### Abstract

We present SOCIAL CHEMISTRY, a new conceptual formalism to study people's everyday social norms and moral judgments over a rich spectrum of real life situations described in natural language. We introduce SOCIAL-CHEM-101, a large-scale corpus that catalogs 292k **rules-of-thumb** such as “It is rude to run a blender at 5am” as the basic conceptual units. Each rule-of-thumb is further broken down with 12 different dimensions of people's judgments, including social judgments of go



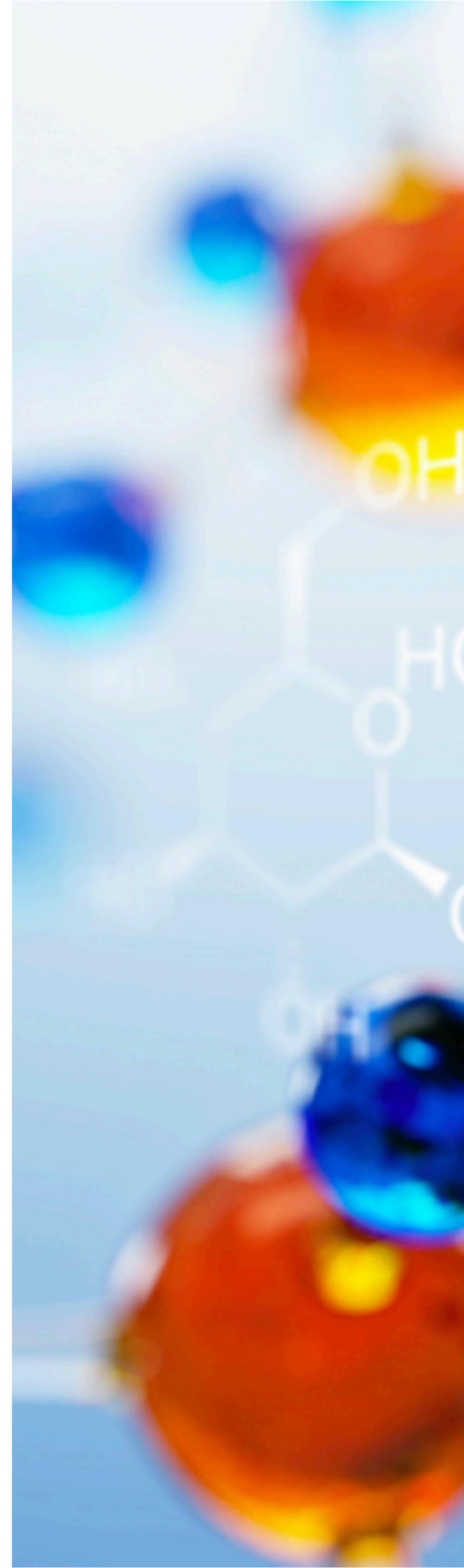
state-of-the-art neural models demonstrate that

Figure illustrates an intuitive subset of  
to reason about social norms in lan-  
guage. Our approach centers around Rules-of-Thumb  
(RoTs: text in colored tubes), which describe social ex-





Narrator: asking my boyfriend to stop being friends with his ex





## Rules of Thumb

It's okay to ask your significant other to stop doing something you're uncomfortable with.

## Situation

Narrator: asking my boyfriend to stop being friends with his ex





RoT

It's not right to tell another person who to spend time with.

Rules of Thumb

It's okay to ask your significant other to stop doing something you're uncomfortable with.

Situation

Narrator: asking my boyfriend to stop being friends with his ex

RoT

You should make sure your SO doesn't feel like a lower priority than your ex.



RoT

It's not right to tell another person who to spend time with.

bad

authority/subversion

Rules of Thumb

It's okay to ask your significant other to stop doing something you're uncomfortable with.

expected/OK

care/harm

Situation

Narrator: asking my boyfriend to stop being friends with his ex

RoT

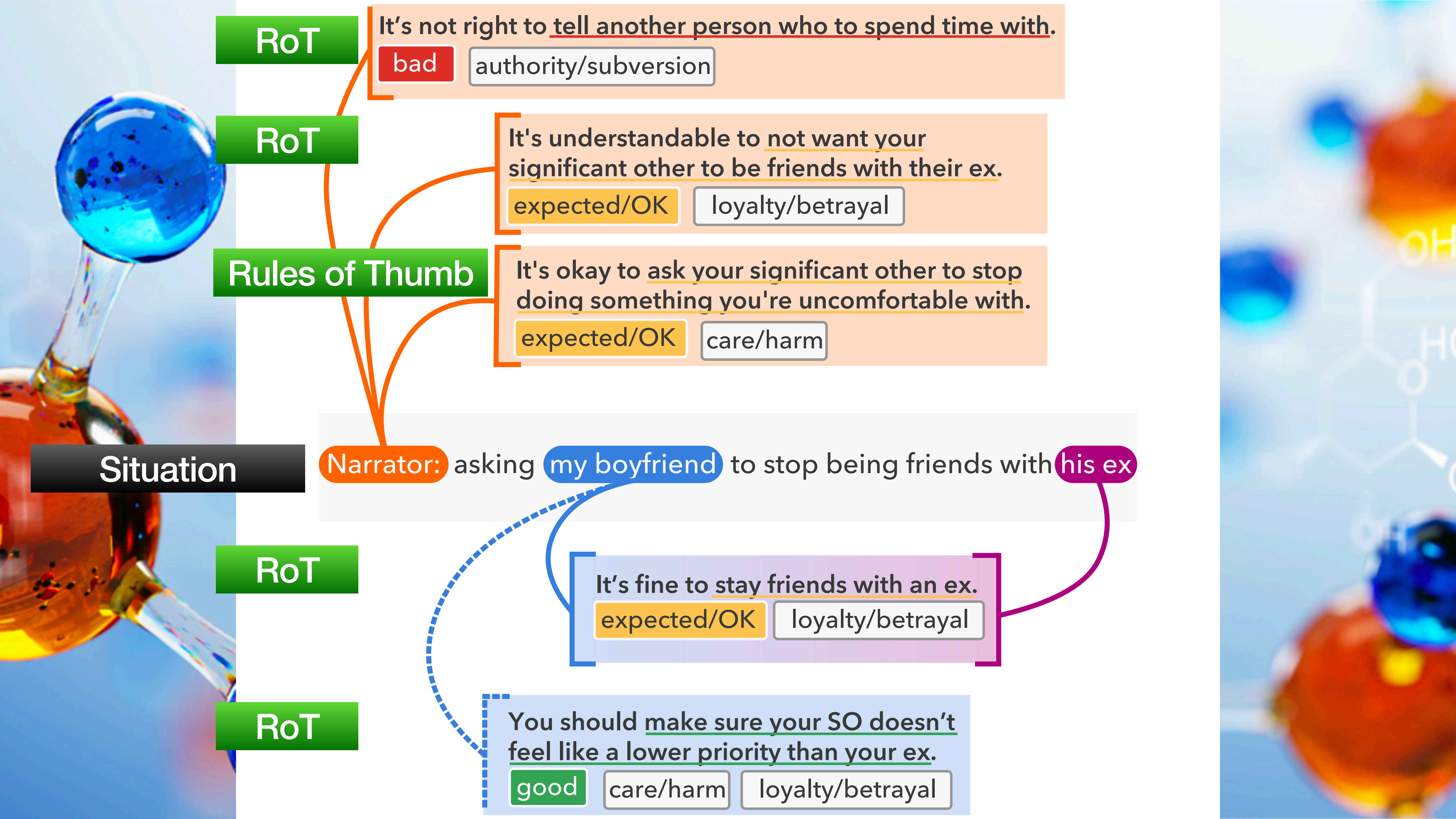
You should make sure your SO doesn't feel like a lower priority than your ex.

good

care/harm

loyalty/betrayal







# Situation

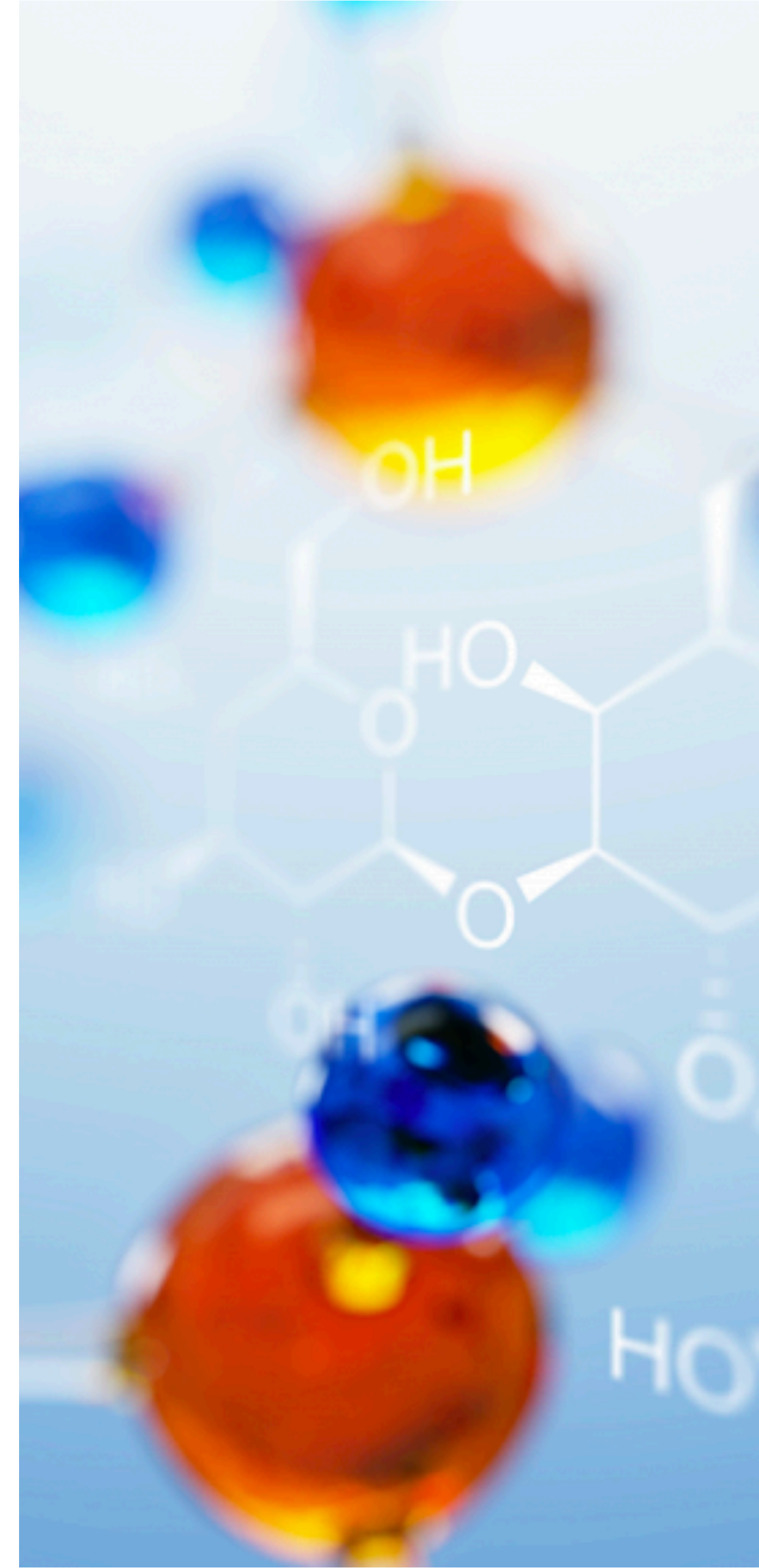
RoT

SITUATION

**Narrator:** Not wanting to be around my GF when she's sick

ROT

It's kind to sacrifice your well-being to take care of a sick person.





Situation

RoT

SITUATION

**Narrator:** Not wanting to be around my GF when she's sick

ROT

It's kind to sacrifice your well-being to take care of a sick person.

ATTRIBUTE KEY

Grounded  
Social

ROT BREAKDOWN

ANTICIPATED AGREEMENT (ROT)

< 1%

~5% - 25%

~ 50%

~ 75% - 90%

> 99%

ROT CATEGORIZATION

Morality / Ethics

Social Norms

Advice

It is what it is

MORAL FOUNDATIONS

Care / Harm

Fairness / Cheating

Loyalty / Betrayal

Authority / Subversion

Sanctity / Degradation

ROT TARGETING

narrator

my GF

no one listed

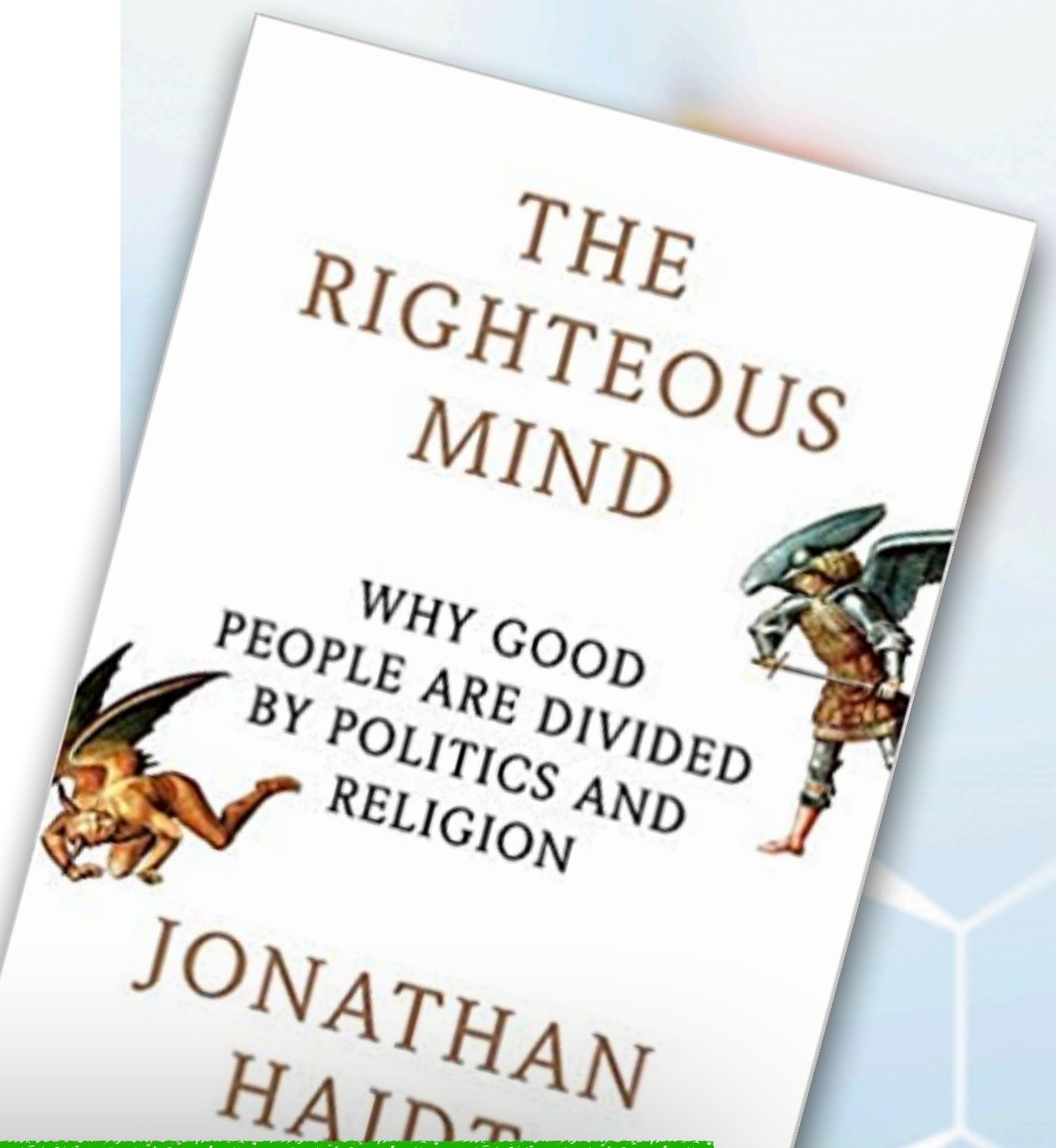
ACTION BREAKDOWN

ACTION

sacrificing your well-being to take care of a sick person

AGENCY

ORIGINAL JUDGMENT



300,000 Rules of Thumb  
grounded on 104k real life situations  
each RoT with 12 structured attributes

e.g.,



Social Judgment



Cultural Pressure



Agency



Anticipated Agreement



Legality



Moral Foundation

Explicitly not

Probably not

Hypothetical

Probable

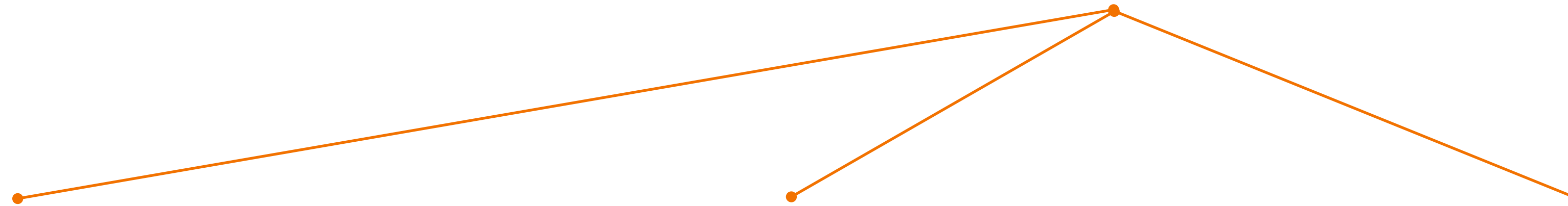
Explicit







# Communication requires understanding **people's norms**.



## **Social Norms**

**E.G.** It's rude to make loud noises during the night.

## **Moral Norms**

**E.G.** It's important to be considerate of others' physical needs, like sleep.

## **Ethical Norms**

**E.G.** You should follow household rules for everyone's benefit.

## **Social Norms**

E.G. It's rude to make loud noises during the night.

## **Moral Norms**

E.G. It's important to be considerate of others' physical needs, like sleep.

## **Ethical Norms**

E.G. You should follow household rules for everyone's benefit.

## **Challenging to capture**

Who decides? Where do we find them?

## **Challenging to utilize**

How do we use abstract moral rules?



# DESCRIPTIVE NORMS



## **Challenging to capture**

Who decides? Where do we find them?

## **Challenging to utilize**

How do we use abstract moral rules?

## **Descriptive norms**

*E.g., descriptive ethics: have everyday people describe*

# GROUNDED DESCRIPTIVE NORMS



## **Challenging to capture**

Who decides? Where do we find them?



## **Challenging to utilize**

How do we use abstract moral rules?

## **Descriptive norms**

*E.g., descriptive ethics: have everyday people describe*

## **Grounded**

*Evoked from situations, tied to people involved*



# In this talk: Reasoning as Generation

- **Part 1:** unsupervised inference-time algorithms

Reasoning thru  
**Neural Backpropagation**

DeLorean

Reasoning thru  
**Search with Logical Constraints**

NeuroLogic

Reasoning thru  
**Distributional Neural Imagination**

Reflective Decoding

- **Part 2:** supervision with declarative knowledge for knowledge modeling

**Neural & Symbolic  
Commonsense Knowledge**

COMET & ATOMIC 2020

**Visually Grounded  
Commonsense Knowledge**

Visual COMET

**Social, Ethical, Moral Norms**

Social Chemistry 101

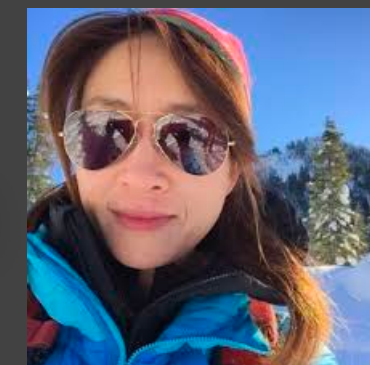
- **Part 3:** benchmarks and algorithmic bias reduction

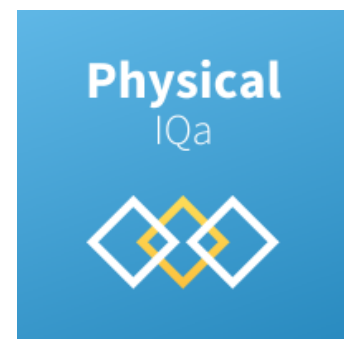


"Got  
a Challenge  
Dataset?"



"Yes!"





1. Physical IQA (**AAAI 2020**)



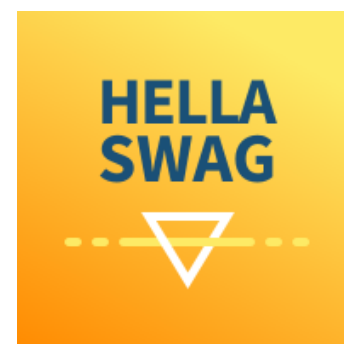
2. Social IQA (**EMNLP 2019**)



3. Visual Commonsense Reasoning (**CVPR 2019**)



4. Abductive Commonsense Reasoning (**ICLR 2020**)



5. HellaSwag (**ACL 2019**)



6. Winogrande (**AAAI 2020**)



7. Cosmos QA (**EMNLP 2019**)

=



A rainbow of  
commonsense  
challenges



# Physical IQA

Test knowledge of affordances and physical attributes

*Q. Which household item can I use to roll out bread dough?*

Wine Bottle

Bread Mixer

~21,000 Instances

GRANDE



473ml

Best Paper Award @ AAAI 2020

# WinoGrande: Adversarial Winograd Schema Challenge at Scale



“The large **ball** crashed right through the **table** because **it** was made of steel.”

“The large **ball** crashed right through the **table** because **it** was made of styrofoam.”



~44,000 Instances



# Cosmos QA

Context:

*It 's a very humbling experience when you need someone to dress you every morning, tie your shoes, and put your hair up. Every menial task takes an unprecedented amount of effort .*

Question:

*What's a possible reason the writer needed someone to dress him every morning?*

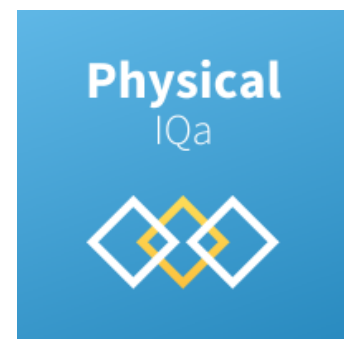
the writer is bad at doing his own hair.

the writer doesn't like putting effort into these tasks.

the writer has a physical disability.

the writer has never learned how to get dressed.

~35,000 Instances



1. Physical IQA (**AAAI 2020**)



2. Social IQA (**EMNLP 2019**)



3. Visual Commonsense Reasoning (**CVPR 2019**)



4. Abductive Commonsense Reasoning (**ICLR 2020**)



5. HellaSwag (**ACL 2019**)



6. Winogrande (**AAAI 2020**)



7. Cosmos QA (**EMNLP 2019**)

=



A rainbow of  
commonsense  
challenges



# The New York Times

## Finally, a Machine That Can Finish Your Sentence

Completing someone else's thought is not an easy trick for A.I. But new systems are starting to crack the code of natural language.



By Cade Metz

Nov. 18, 2018

In August, researchers from the Allen Institute for Artificial Intelligence, a lab based in Seattle, unveiled a new system called [BART](#). It examined whether machines could generate text like this one:

On stage, a woman takes a seat at the piano. She



What's going on???



Sebastian Ruder  
@seb\_ruder

Following

It's amazing how fast [#NLProc](#) is moving these days. We have now reached super-human performance on SWAG, a commonsense task that will only be introduced at @emnlp2018 in November! We need even more challenging tasks! BERT: [arxiv.org/abs/1810.04805](#) SWAG: [arxiv.org/abs/1808.05326](#)

System	Dev	Test
ESIM+GloVe	51.9	52.7
ESIM+ELMo	59.1	59.2
BERT <sub>BASE</sub>	81.6	-
BERT <sub>LARGE</sub>	<b>86.6</b>	<b>86.3</b>
Human (expert) <sup>†</sup>	-	85.0
Human (5 annotations) <sup>†</sup>	-	88.0

Table 4: SWAG Dev and Test accuracies. Test results are scored against the hidden labels by the SWAG model. <sup>†</sup>Human performance is measure with 100%, as reported in the SWAG paper.

On stage, a woman takes a seat at the piano. She  
a) sits on a bench as her sister plays with the doll.  
b) smiles with someone as the music plays.  
c) is in the crowd, watching the dancers.  
d) **nervously sets her fingers on the keys.**  
A girl is going across a set of monkey bars. She  
a) jumps up across the monkey bars.  
b) struggles onto the monkey bars to grab her head.  
c) **gets to the end and stands on a wooden plank.**  
d) jumps up and does a back flip.  
The woman is now blow drying the dog. The dog  
a) **is placed in the kennel next to a woman's feet.**  
b) washes her face with the shampoo.  
c) walks into frame and walks towards the dog.  
d) tried to cut her face, so she is trying to do something very close to her face.

Table 1: Examples from *SWAG*; the correct answer is **bolded**. Adversarial Filtering ensures that stylistic models find all options equally appealing.

5:38 AM - 12 Oct 2018

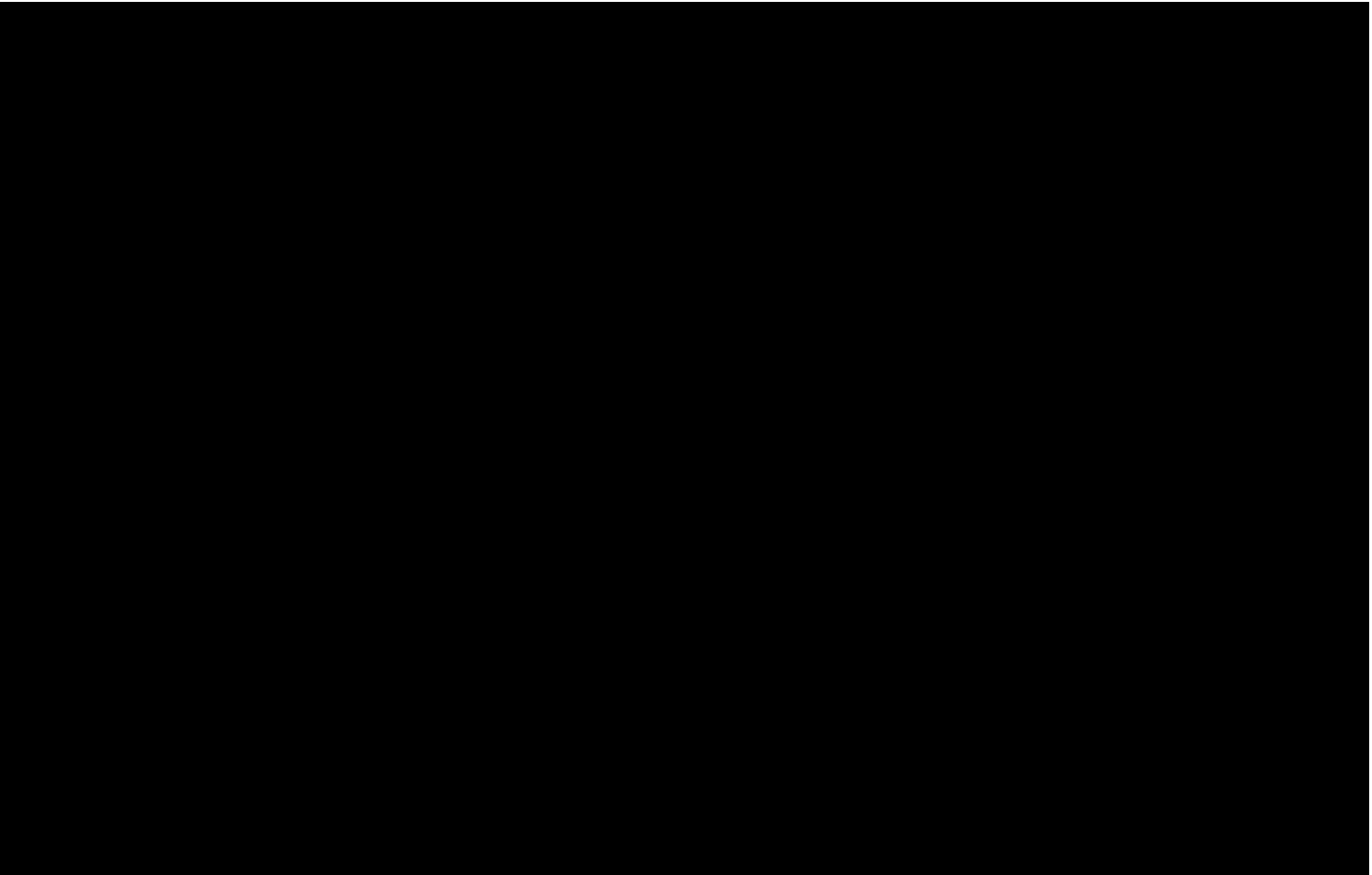


Thomas Wolf  
@Thorn\_Wolf

Following

BERT is super impressive! Amazing development of the nice OpenAI GPT!

Human level already reached on the recent SWAG dataset (EMNLP'18)! I'm wondering if we should consider the task "solved" or if we could/should update such an adversarially generated dataset?





Many AI systems  
perform well,  
but do so for  
**questionable**  
**reasons**





# HellaSwag: Can a Machine Really Finish Your Sentence?

ACL 2019

- What happened with SWAG? TLDR:
- dataset must be debiased (by algorithms)
  - dataset must evolve (with the evolving SOTA)

Rowan  
Zellers



Ari  
Holtzman



Yonatan  
Bisk



Ali  
Farhadi



Yejin Choi



*Train  $t_1$*

*Test  $t_1$*

**Correct  
Answer**

I

**Potential  
Incorrect  
Answers**

I

$x_1^+$

$x_2^+$

$x_3^+$

$x_4^+$



$x_5^+$



$x_6^+$

• • •

$x_{1,1}^-$

$x_{4,1}^-$



$x_{5,1}^-$

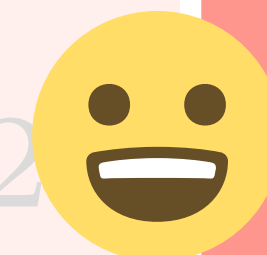


$x_{6,1}^-$

• • •

$x_{1,2}^-$

$x_{4,2}^-$



$x_{5,2}^-$



$x_{6,2}^-$

• • •

$x_{1,3}^-$

$x_{2,3}^-$

$x_{3,3}^-$

$x_{4,3}^-$



$x_{5,3}^-$



$x_{6,3}^-$

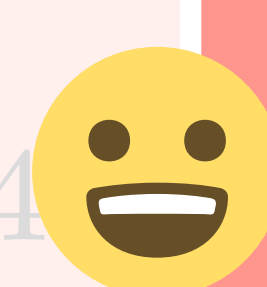
• • •

$x_{1,4}^-$

$x_{2,4}^-$

$x_{3,4}^-$

$x_{4,4}^-$



$x_{5,4}^-$



$x_{6,4}^-$

• • •

•  
•  
•

•  
•  
•

•  
•  
•

•  
•  
•

•  
•  
•

•  
•  
•

(the filtering model)





*Test  $t_2$*

*Train  $t_2$*

**Correct  
Answer**

**I**

$x_1^+$

$x_2^+$

$x_3^+$

$x_4^+$

$x_5^+$

$x_6^+$

$x_{1,1}^-$

$x_{2,1}^-$

$x_{3,1}^-$

$x_{6,1}^-$

$x_{1,2}^-$

$x_{2,2}^-$

$x_{3,2}^-$

New  $x_{6,2}^-$

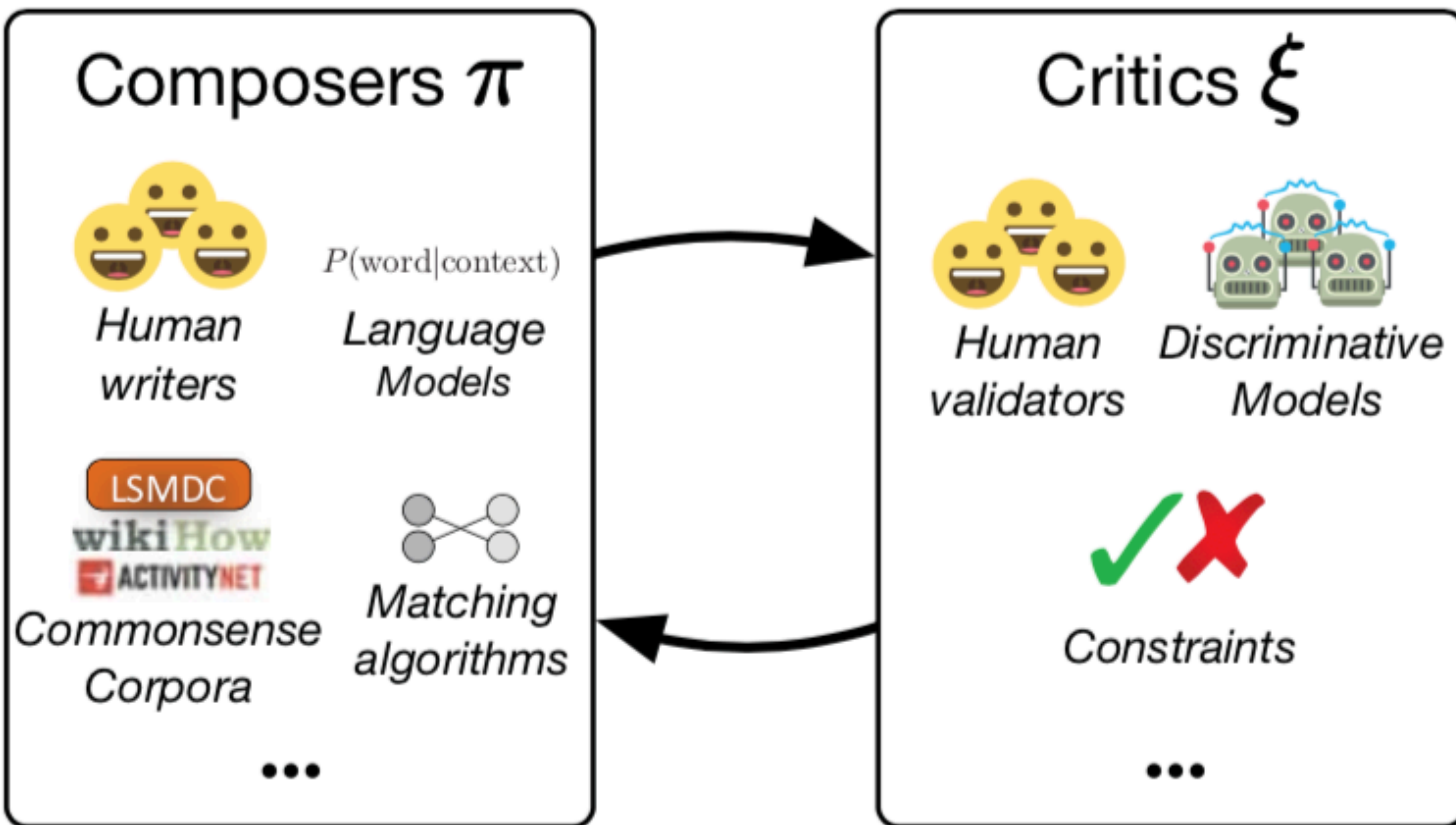
**Potential  
Incorrect**

**T**



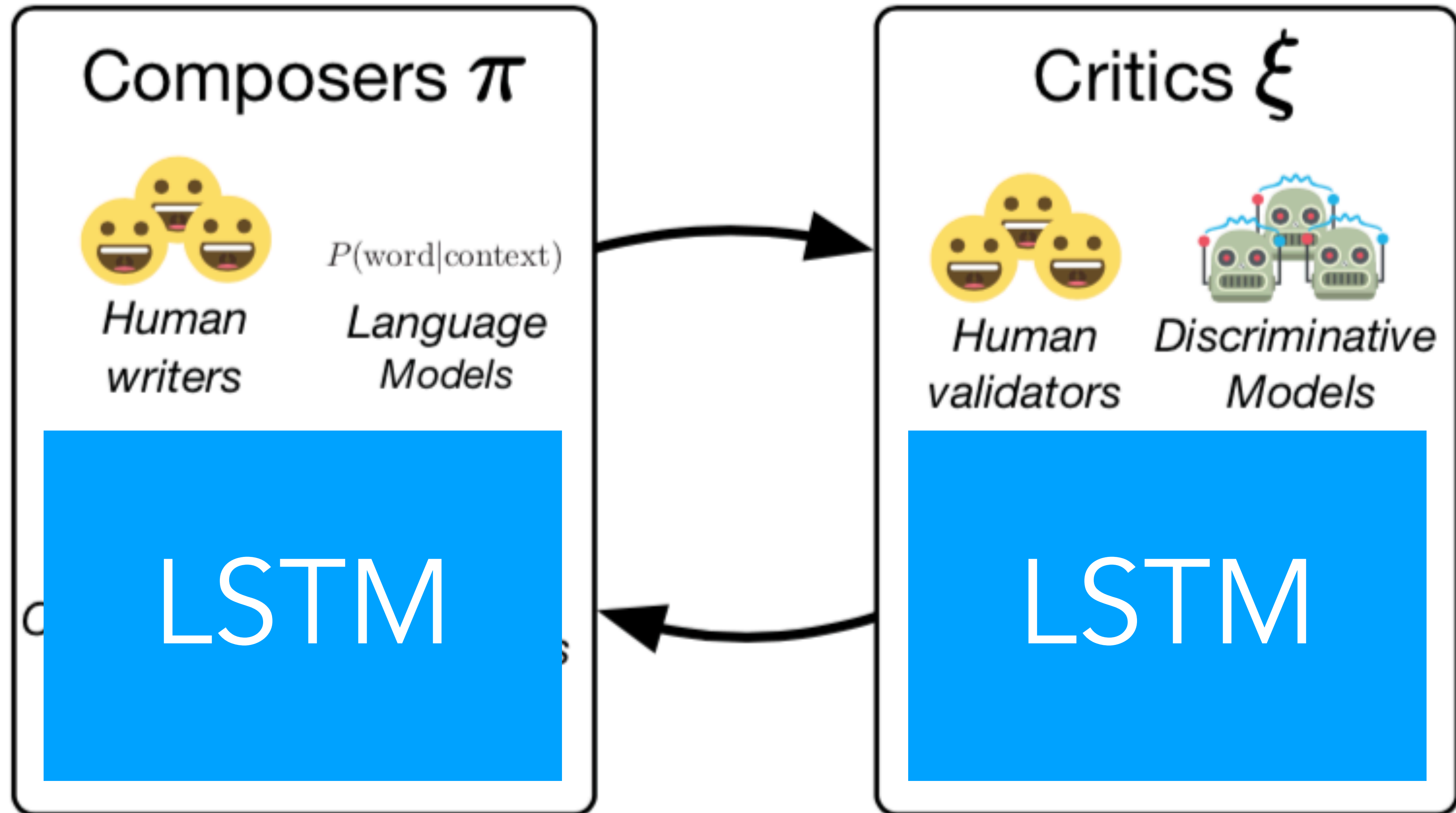
- We repeat this until convergence!
- The resulting AF'ed data holds identical distribution across training & test
- Human performance remains high on the AF'ed data, while AI performance drops considerably

# An Adversarial Evaluation Framework of

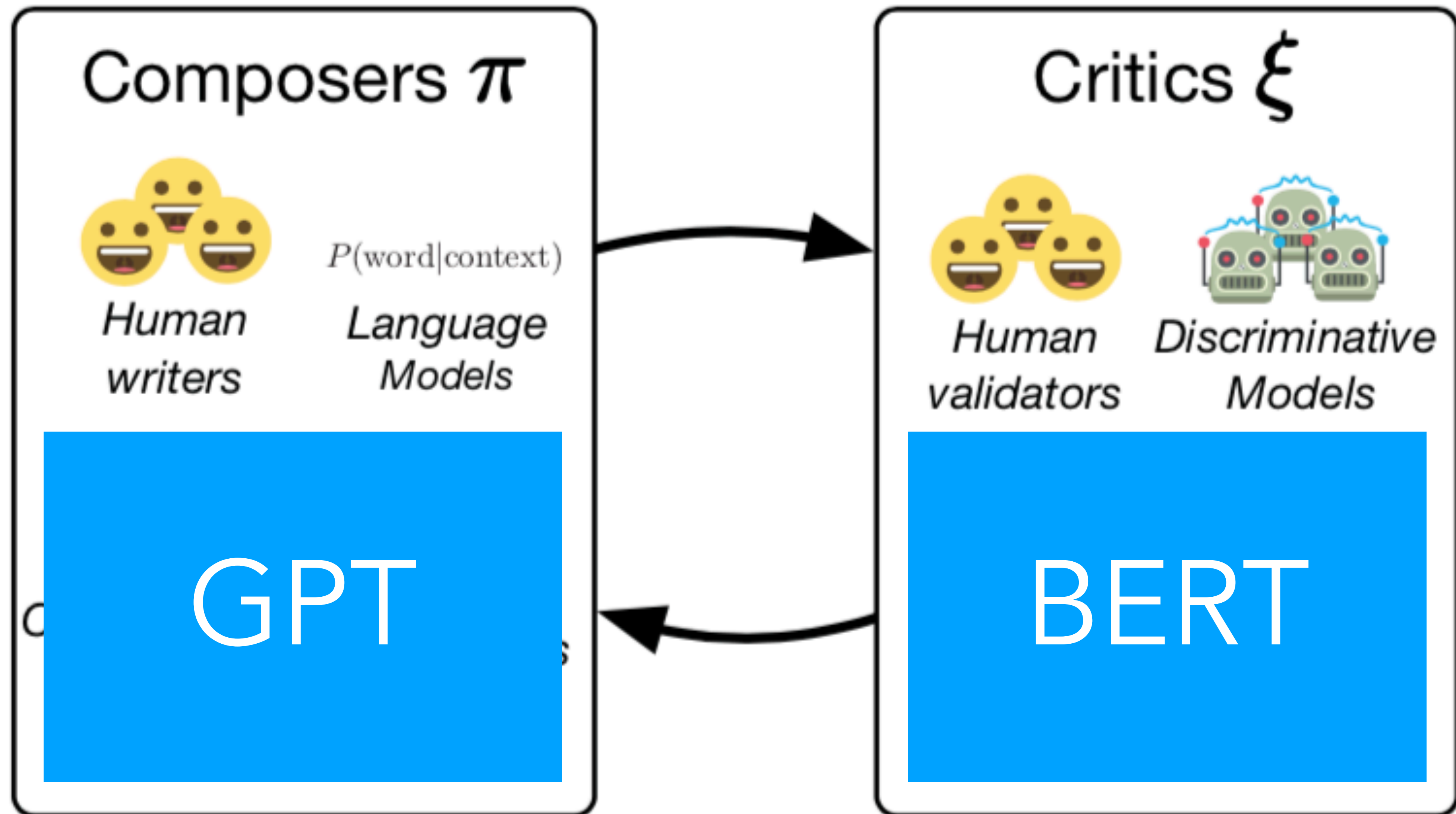




# Pre GPT and BERT era (swag, 2018)

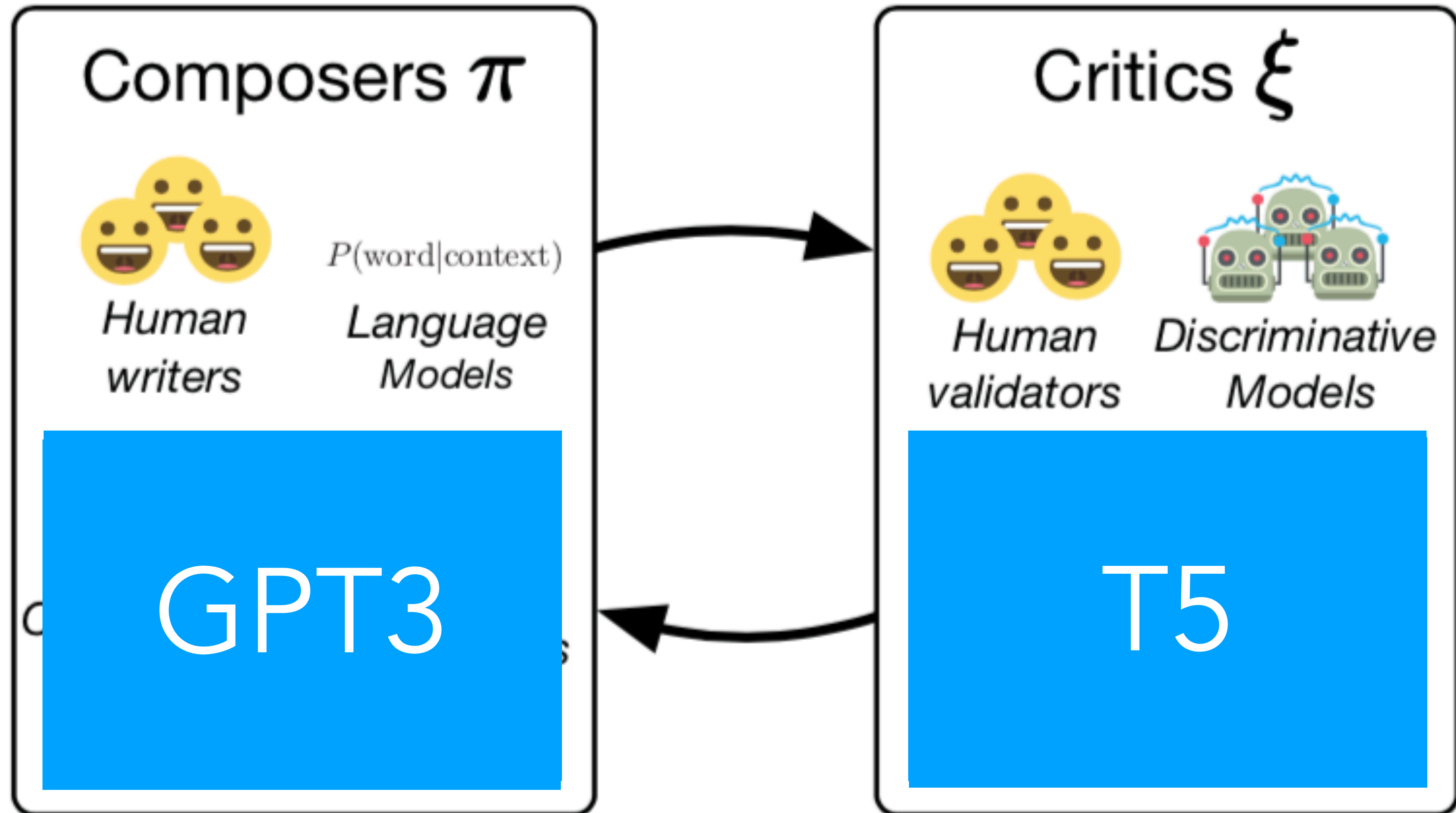


# Post GPT and BERT era (hellaswag 2019)





# Post GPT and BERT era (2020?)



How do we check if a model only solved a “dataset”  
without solving an underlying task? Try adversarial evaluation!

data

AF

benchmarks must  
evolve!

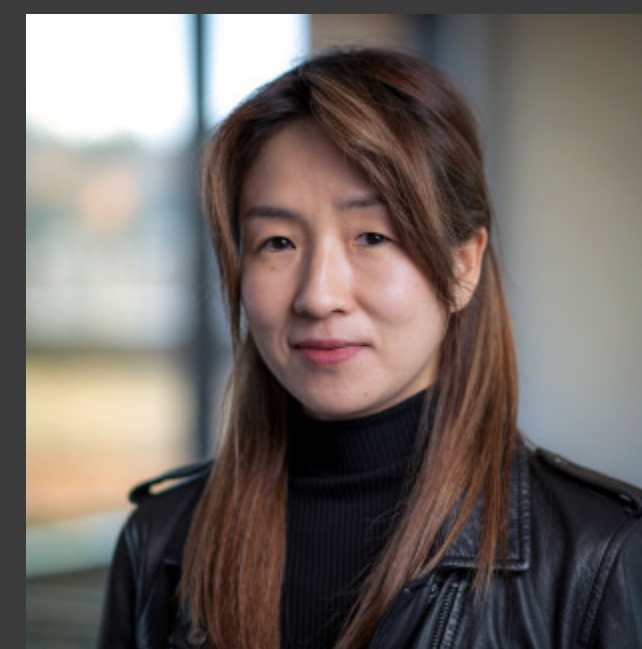




# WinoGrande

Adversarial Winograd Schema Challenge at Scale

**Outstanding Paper (Best Paper) Award at AAAI'20**







# Adversarial Filters of Dataset Biases

ICML 2020



Ronan LeBras



Swabha Swayamdipta



Chandra  
Bhagavatula



Rowan Zellers



Matt  
Peters



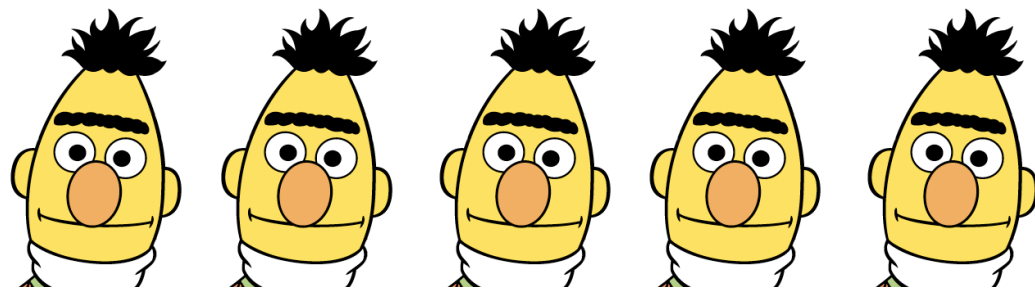
Ashish  
Sabharwal



Yejin  
Choi



Original  
Flavor!



- **SWAG** (Zellers et al., 2018)
- **HellaSWAG** (Zellers et al., 2019)

Af-Optimum  
Flavor!



Af-Lite  
Flavor!

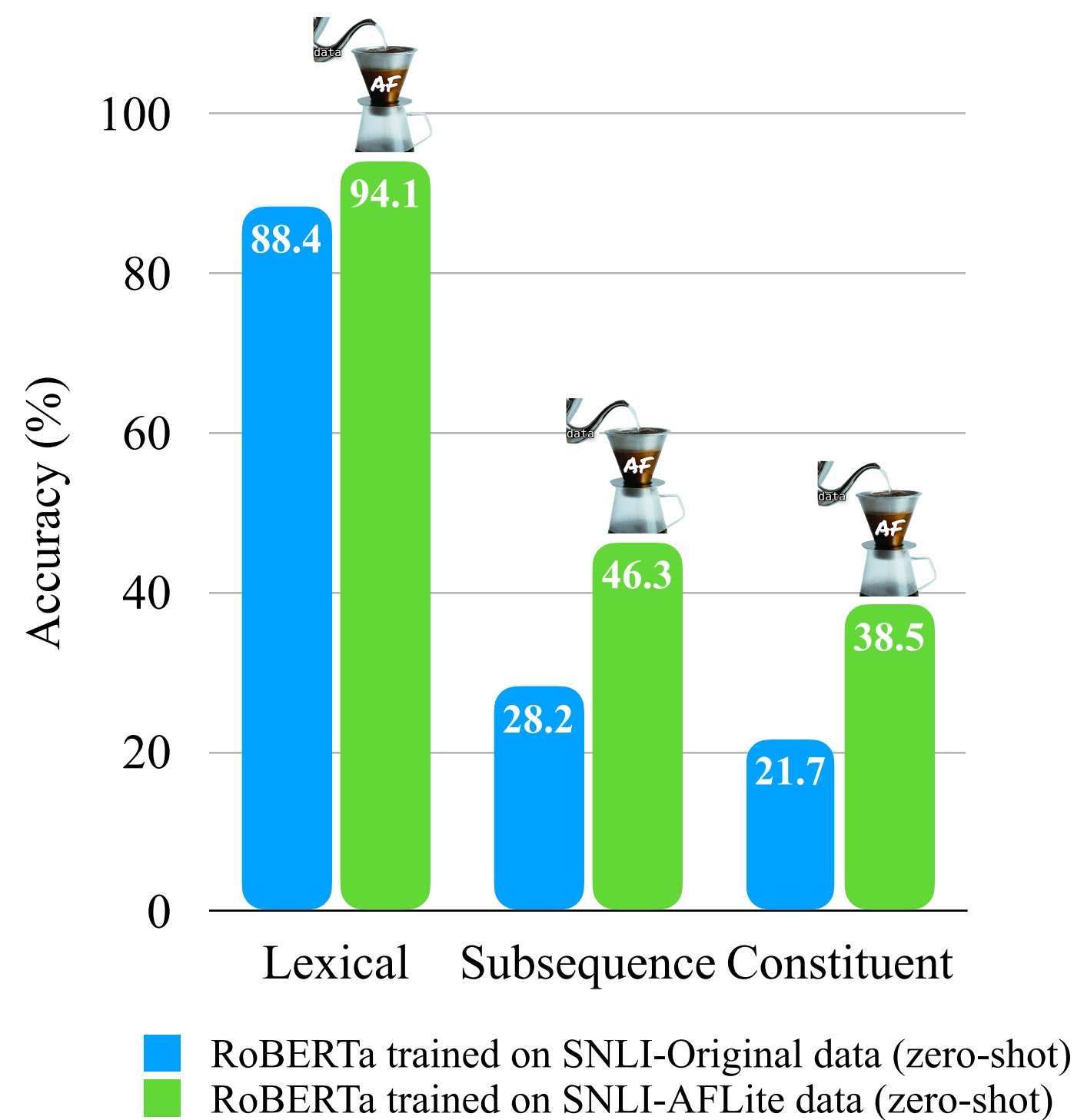


- **Winogrande** (Sakaguchi et al., AAAI 2020)
- **Adversarial Filters of Dataset Biases**  
(Le Bras et al., ICML 2020)

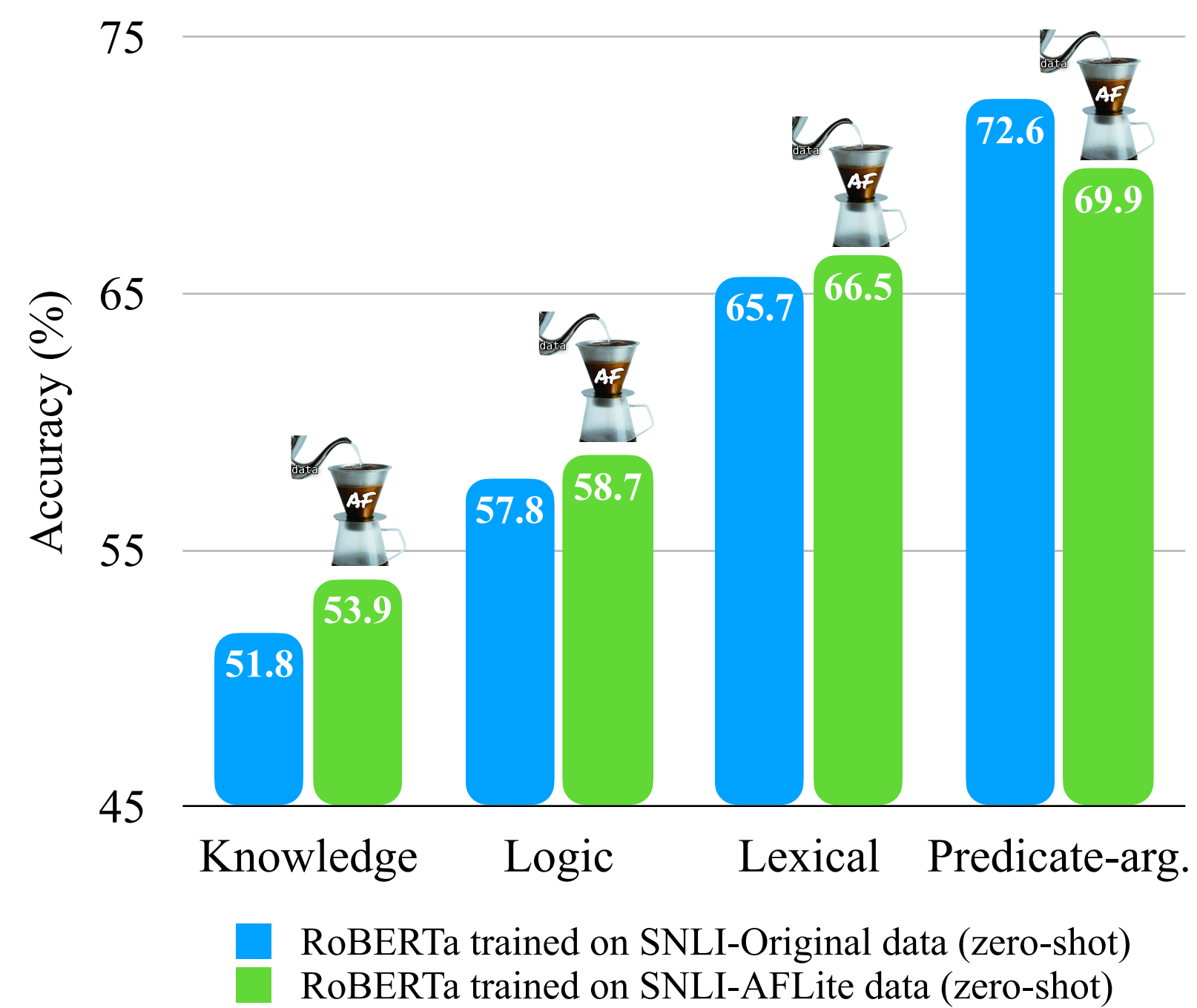


# Models Trained on AF'ed Data — Better on Out-of-Distribution

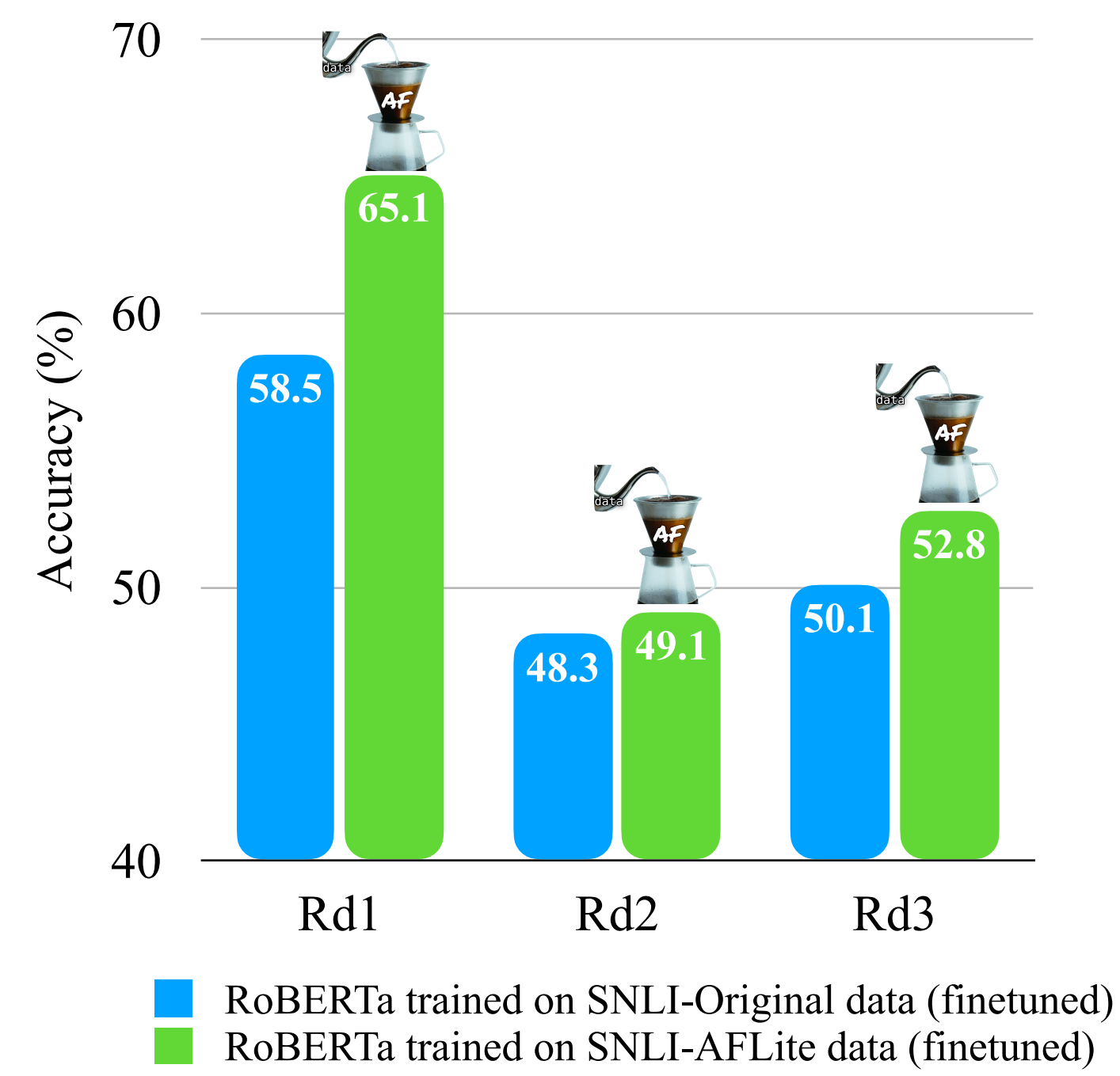
## HANS



## NLI-Diagnostics



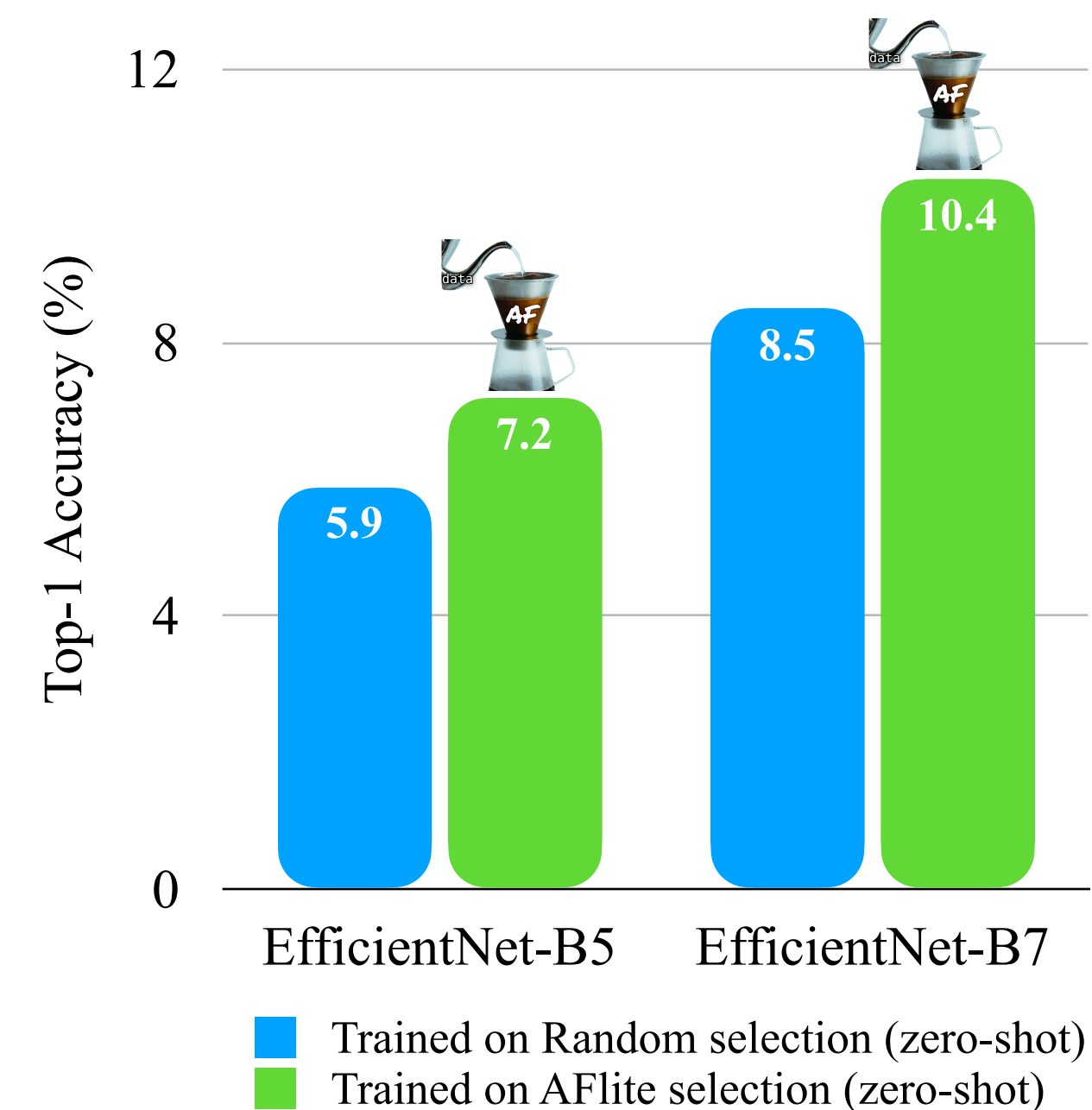
## Adversarial-NLI





# Models Trained on AF'ed Data — Better on Out-of-Distribution

ImageNet-A (Hendrycks et al., arXiv'19)





# Current Paradigm of QA datasets

1. Neural networks immediately latch on to spurious correlations & unwanted dataset biases
  - Multiple-choice QA datasets might require algorithmic debiasing to avoid overestimation of true AI capabilities
2. Can you imagine learning only by solving a lot of exam problems? — it's hard even for humans!
  - Need to **evaluate** machines **generatively**
  - Need to teach machines **concepts** more **directly**

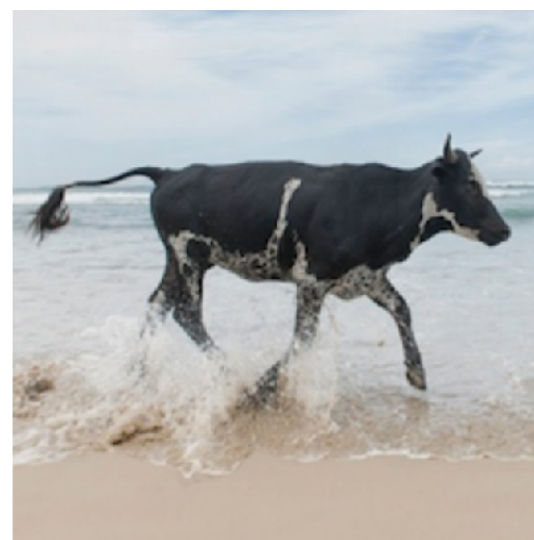
(Beery et al., 2018)



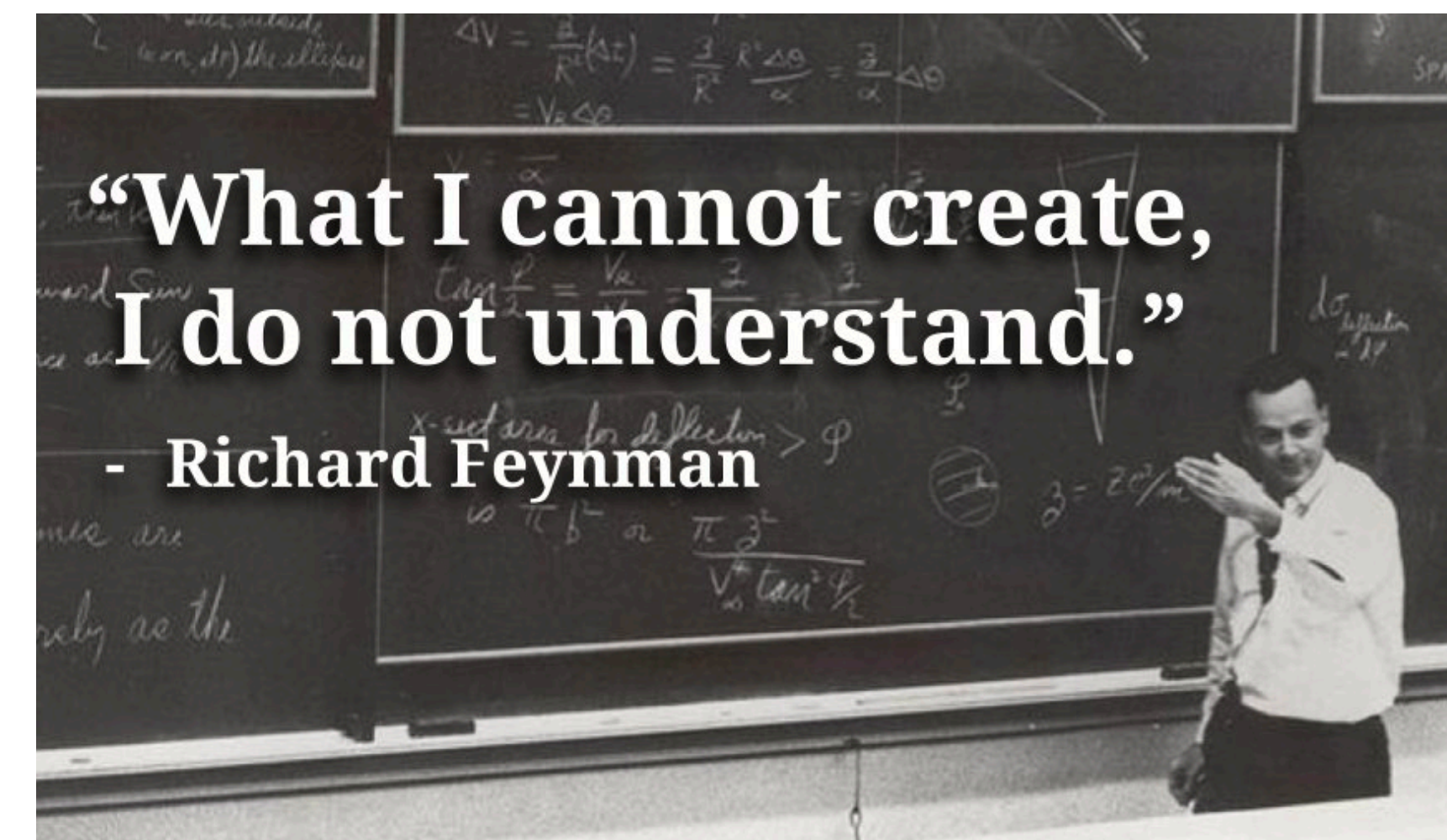
Cow 🤖



Cow 🤖



No cow 🤖



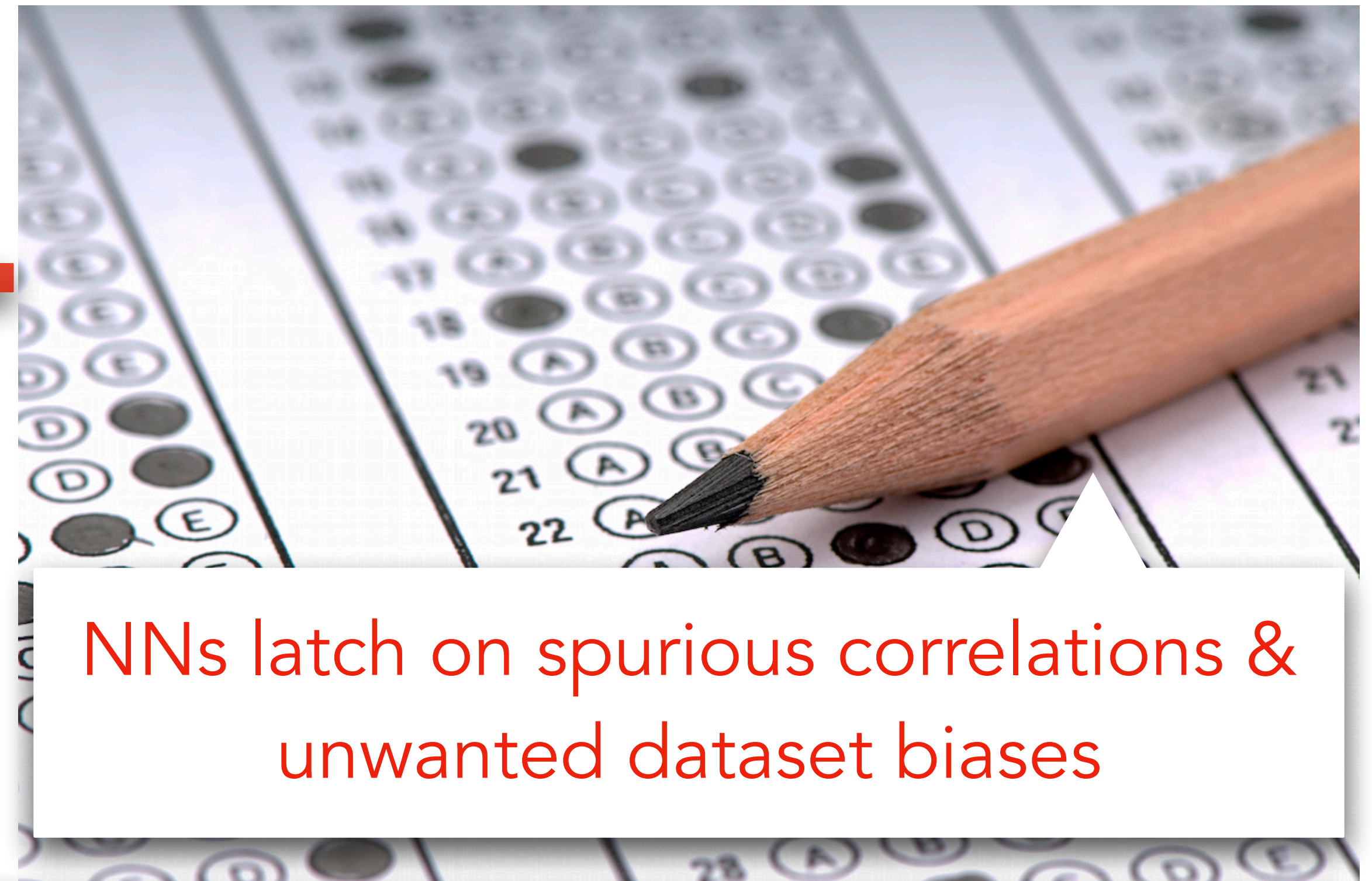
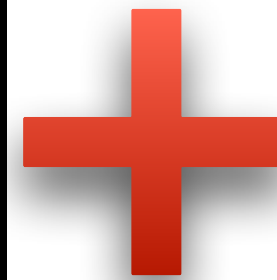


# Imagine taking a deep learning class in which...

**Self-supervision** on a lot of DL code

**Supervision** on a lot of exam problems

```
86 self.names = names
87 self.name2index = dict(zip(names, range(len(names))))
88
89
90 def __del__(self):
91     # free memory created by C to avoid memory leak
92     if hasattr(self, '__createfrom__') and self.__createfrom__ == 'C':
93         if pointer(self) is not None:
94             libbigfile.free_file(pointer(self))
95
96 def read(self, requested, isname=True):
97     if isname:
98         index_name_array = [(self.name2index[x], x) for x in requested]
99     else:
100         assert(min(requested)>=0)
101         assert(max(requested)<len(self.names))
102         index_name_array = [(x, self.names[x]) for x in requested]
103         index_name_array.sort()
104
105     npoints = len(index_name_array)
106     c_index = (c_ulonglong * npoints)()
107     for i in range(npoints):
108         c_index[i] = index_name_array[i][0]
109
110     size = self.ndims * npoints
111     pdata = (c_float * size)()
112     res = libbigfile.seq_read_memory(self, npoints, c_index, pdata)
113     assert(res)
```



NNs latch on spurious correlations & unwanted dataset biases

Can't learn concepts well enough; need to learn from declarative knowledge



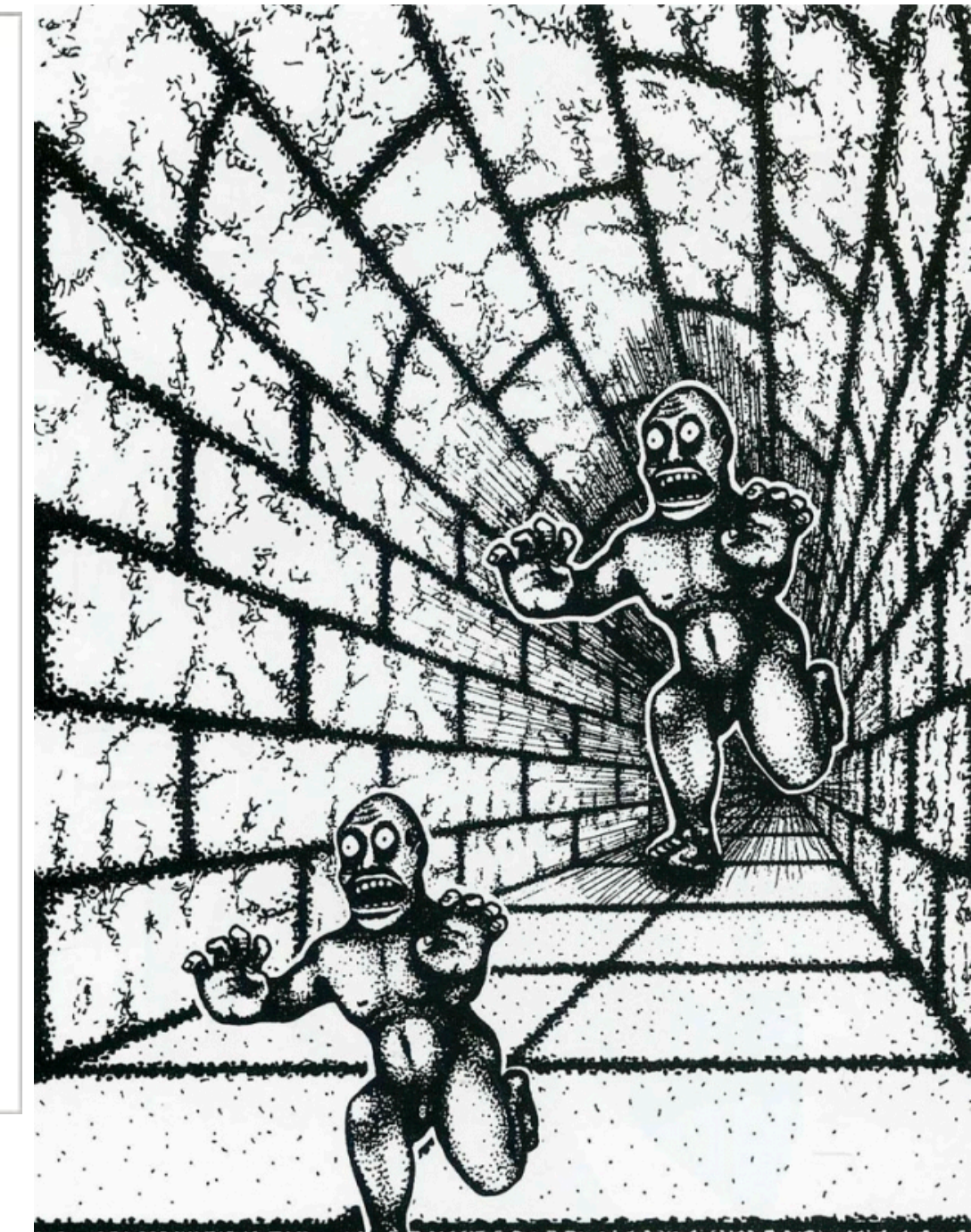
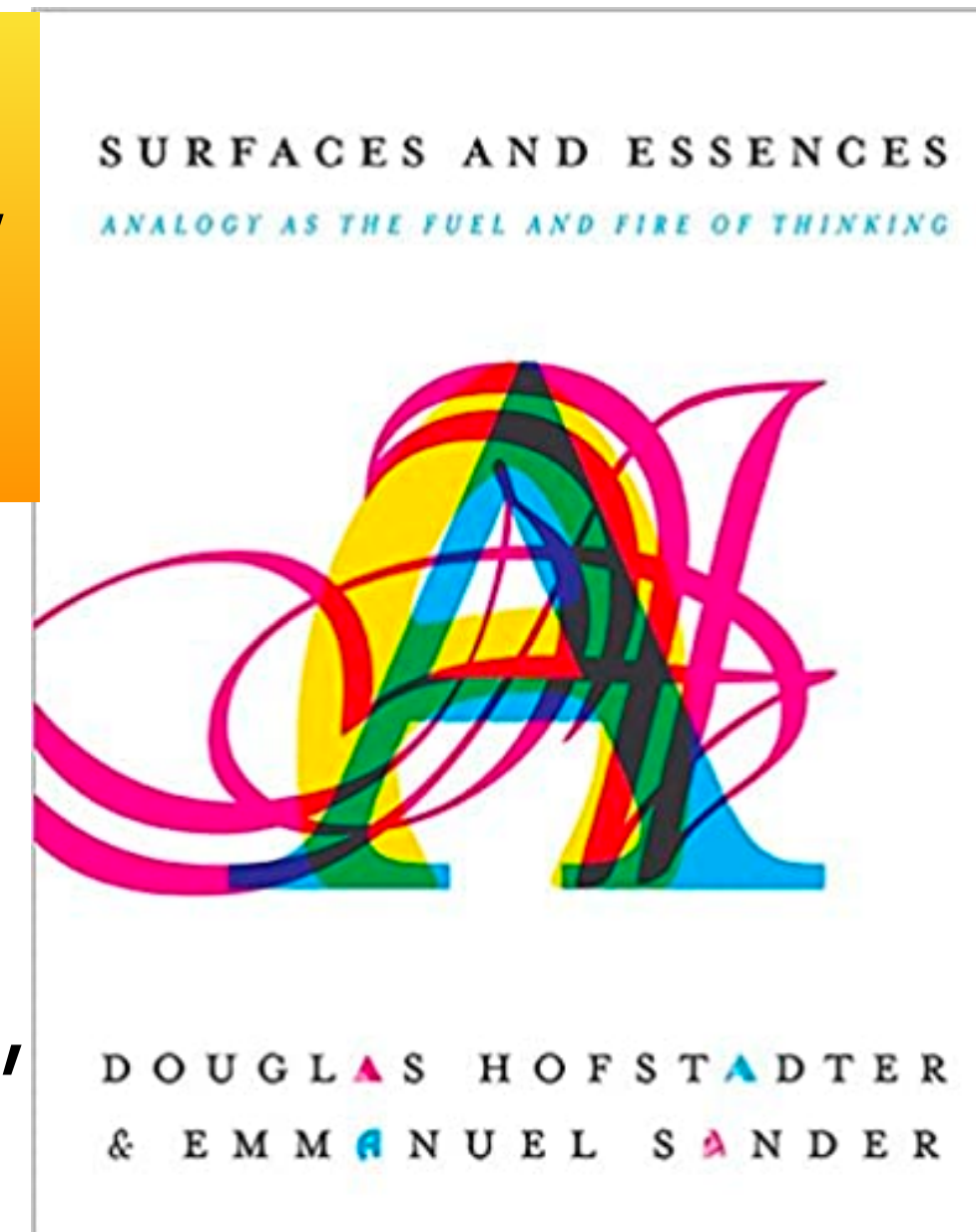
# Tldr; 🔥 Language and Symbols 🔥

## 🔥 Reasoning as Generation 🔥

### 1. *Language* as the *symbols*

- all of it —
- not just **words** (ImageNet labels, ...)
- not just small **sets of words** (scene graphs,

"Categories (concepts) vastly outnumber words, and require free-form open text descriptions"



### 2. Reasoning as *generative tasks*

- As opposed to *discriminative tasks* (i.e., categorization)
- Because the space of reasoning in language is **infinite**

"thinking out loud"

We often think as we speak,  
on the fly, word-by-word  
without enumerating all possible  
alternative sentences





# Commonsense AI: Closing Remarks & Open Research Questions

# Commonsense

- Searching “commonsense” from ACL anthology
  - Most papers are either from 80s or from the past few years

## Position Paper on Common-sense and Formal Semantics

Geoffrey Nunberg  
Xerox PARC and CSLI, Stanford

### 1. A philological excursus

I’m not sure what I’m doing on this panel, but I thought it would be helpful if we could start at the beginning. It’s interesting to note that both the dictionary and common sense were eighteenth-century inventions. This is no coincidence; in fact, it’s entirely appropriate that the most celebrated



# Revisiting Commonsense

I was told not to speak the word commonsense...

Past failures (in 70s – 80s) are inconclusive

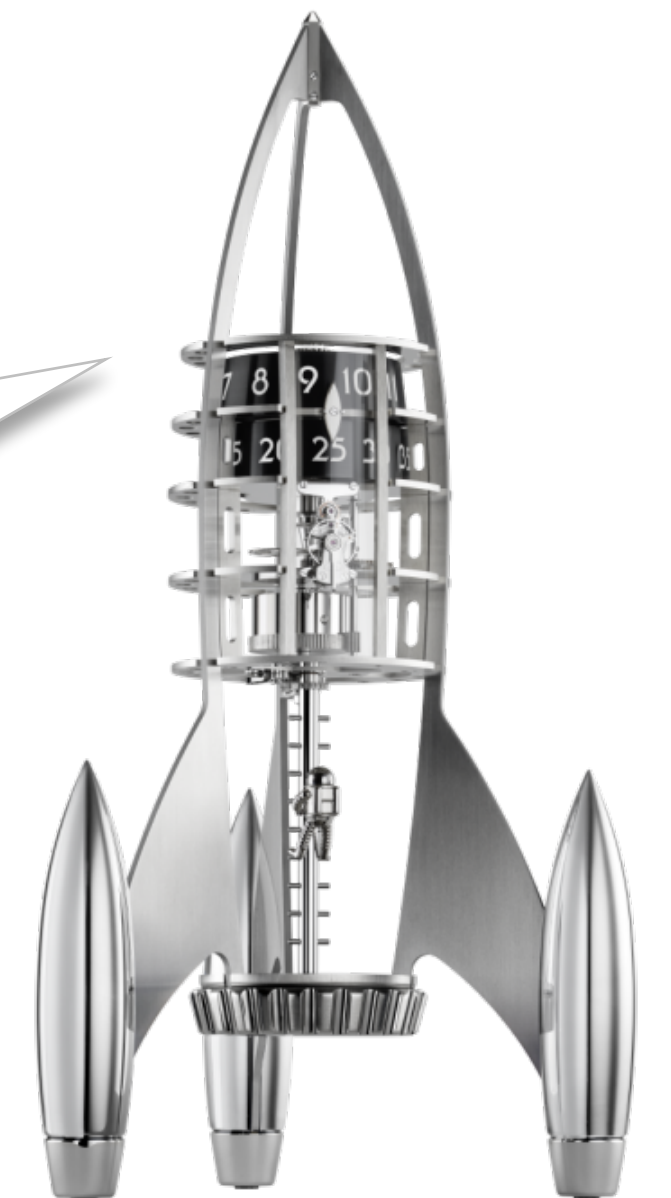
- weak computing power
- not much data
- no crowdsourcing
- not as strong computational models
- not ideal conceptualization / representations

# Path to commonsense?

Brute force larger networks with deeper layers?



You don't reach to the moon  
by making the tallest building in the world  
taller





# Thoughts on Language & Embodiment

- My dog (my cat, my baby...) has common sense without language. Therefore, AI doesn't need language for common sense
  - AI for (ultimately) humans, not for (just) dogs
- Some concepts (e.g., that juicy crunch feel you get when biting into an apple) can't be learned without embodiment. Thus, can't learn concepts without embodied experiences
  - We learn that a tiger can eat a human without experiencing it in real life





# Thoughts on Language & Embodiment

- We might never have a robot that can bite into an apple.
- Nor might we ever have a simulation environment in which AI can experience what it is like for humans to bite on an apple.
- Trade-offs between **the coverage of concepts** and **the richness of multimodal / embodied experiences**
  - 3D environments, real or simulation, allow for richer experiences with a narrow slice of concepts
  - Language, images, videos allow for a significantly broader range of concepts with impoverished embodied experiences
  - **Language** provides a powerful representation medium to learn humans's embodied experiences as **declarative knowledge**





# In this talk: Reasoning as Generation

- **Part 1:** unsupervised inference-time algorithms

Reasoning thru  
**Neural Backpropagation**

DeLorean

Reasoning thru  
**Search with Logical Constraints**

NeuroLogic

Reasoning thru  
**Distributional Neural Imagination**

Reflective Decoding

- **Part 2:** supervision with declarative knowledge for knowledge modeling

**Neural & Symbolic  
Commonsense Knowledge**

COMET & ATOMIC 2020

**Visually Grounded  
Commonsense Knowledge**  
data

Visual COMET

**Social, Ethical, Moral Norms**

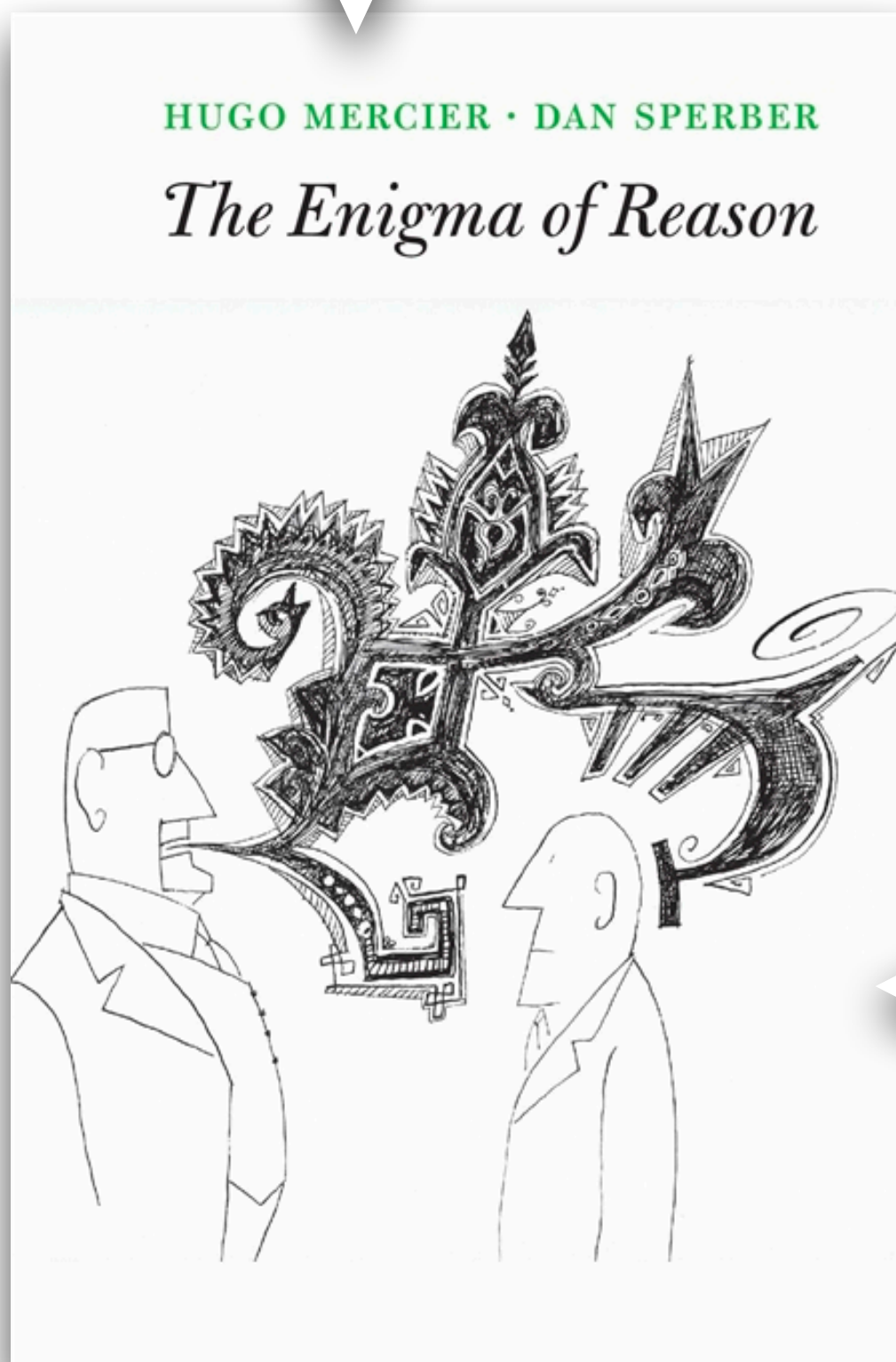
Social Chemistry 101

- **Part 3:** benchmarks and algorithmic bias reduction



# Remarks on Language & Reason

**Intuitive inferences** is a great deal about extraction of new information from the information already available.

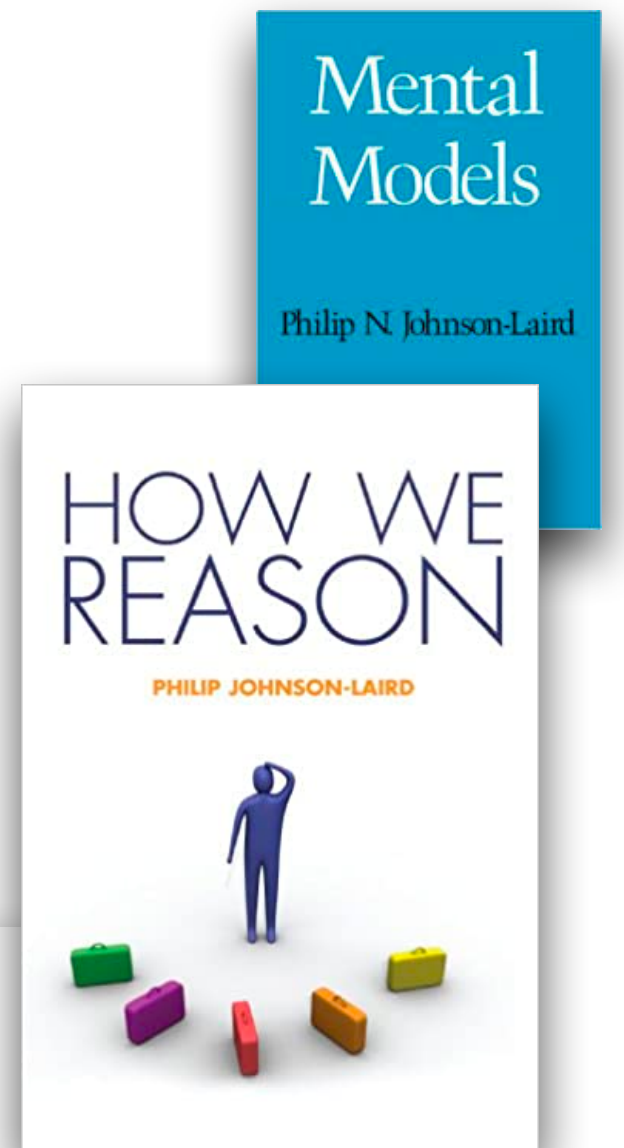


INTUITION  
SYSTEM 1

REASONING  
SYSTEM 2

**Reasons** are used primarily not to guide oneself but to justify oneself in the eyes of others, and to convince others

**Reasoning** serves the purpose of communication



**Human reason** is a mechanism of **intuitive inferences**

... in which logic plays at best a marginal role.



# ACL 2020 Commonsense Tutorial

— 2nd most popular (among 8 tutorials) —

<https://homes.cs.washington.edu/~msap/acl2020-commonsense/>

1288 registrations for our tutorial

2819 view counts on our recorded lectures

(Total registration of the main conference = 4972)



Vered  
Shwartz



Maarten  
Sap



Antoine  
Bosselut



Dan Roth



Yejin  
Choi



Thanks! Questions?