# Commonsense Reasoning in Natural Language Processing

## Vered Shwartz

**Guest Lecture, Deep Learning for NLP**

# The Deep Learning Revolution

# The Deep Learning Revolution

**Translation**

Google's AI translation system is approaching human-level accuracy

# The Deep Learning Revolution

**Translation**

Google's AI translation system is approaching human-level accuracy

**Reading Comprehension**

ALIBABA AI BEATS HUMANS IN READING–COMPREHENSION TEST

CHRISTINE CHOU | JULY 9, 2019

# The Deep Learning Revolution

**Translation**

Google's AI translation system is approaching human-level accuracy

**Reading Comprehension**

ALIBABA AI BEATS HUMANS COMPREHENSION TEST

CHRISTINE CHOU | JULY 9, 2019

**Chatbots**

Artificial intelligence / Voice assistants

Your next doctor's appointment might be with an AI

A new wave of chatbots are replacing physicians and providing frontline medical advice—but are they as good as the real thing?

by **Will Douglas Heaven**                    October 16, 2018

# The Deep Learning Revolution

# The Deep Learning Revolution

Does this mean language understanding is nearly solved?

**Translation**

Google's AI translation system is approaching human-level accuracy

**Reading Comprehension**

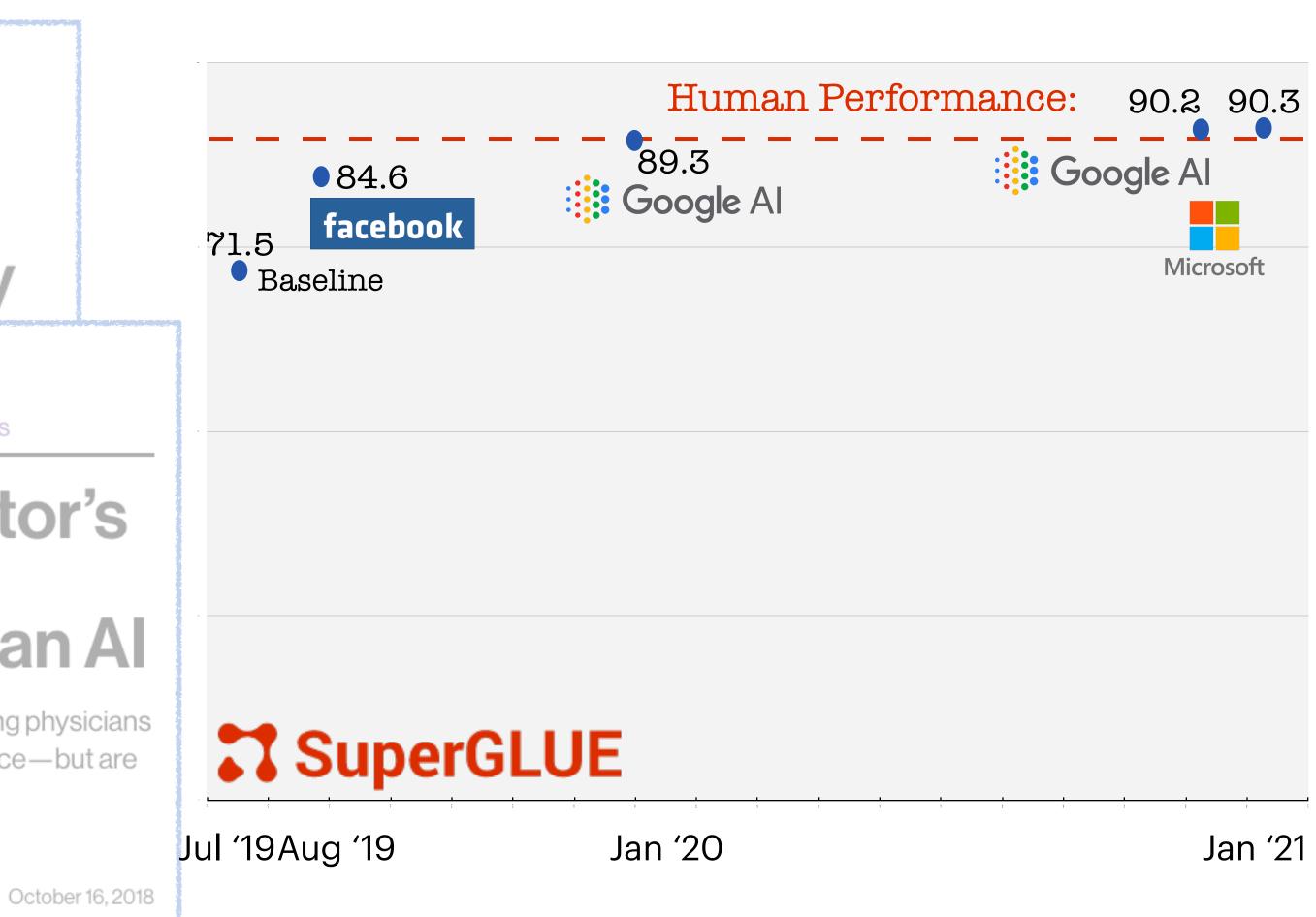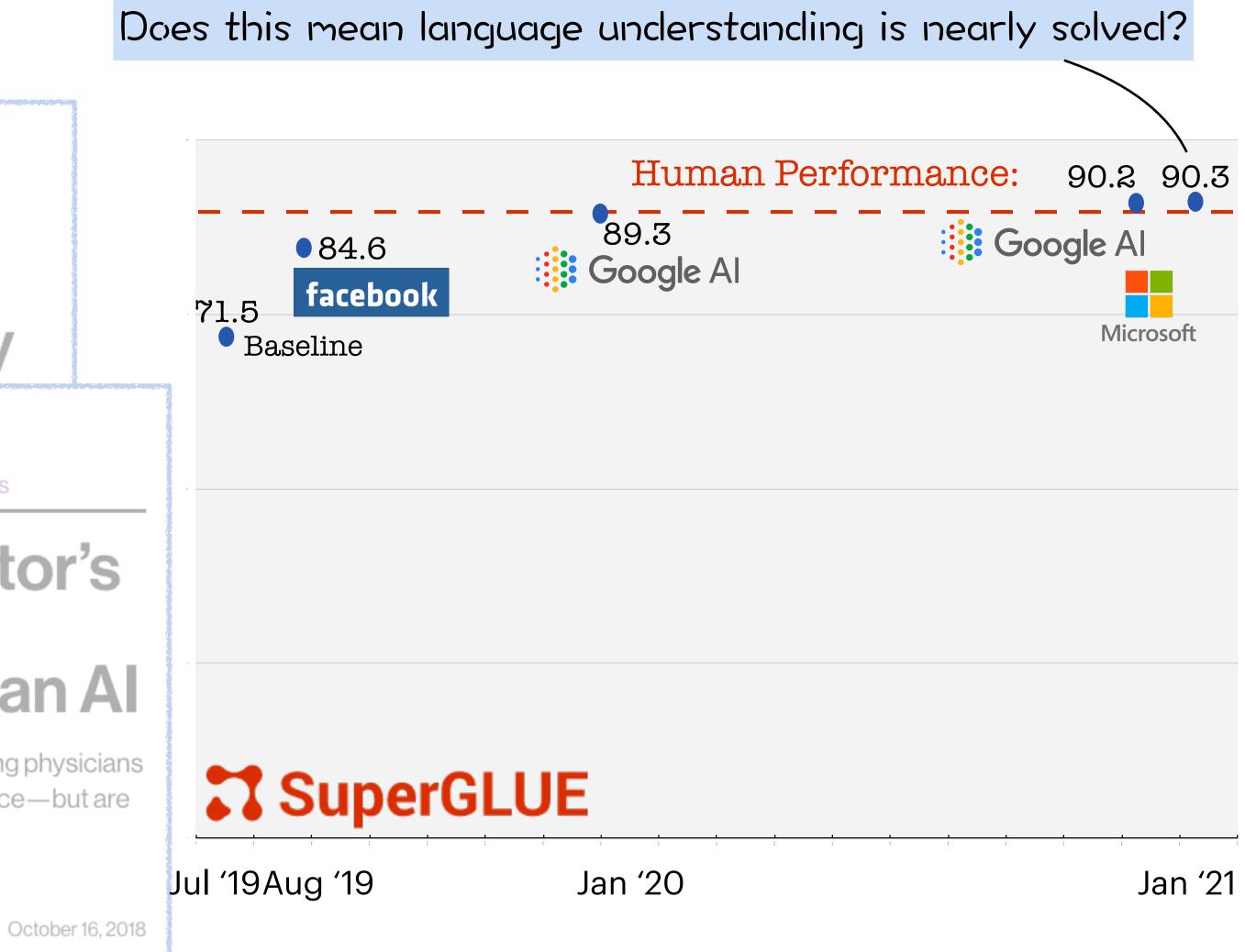ALIBABA AI BEATS HUMANS COMPREHENSION TEST

CHRISTINE CHOU | JULY 9, 2019

**Chatbots**

Artificial intelligence / Voice assistants

Your next doctor's appointment might be with an AI

A new wave of chatbots are replacing physicians and providing frontline medical advice—but are they as good as the real thing?

by **Will Douglas Heaven**     October 16, 2018

Human Performance:     90.2    90.3

84.6    89.3    Google AI

facebook    Google AI

71.5    Microsoft
Baseline

SuperGLUE

Jul '19  Aug '19        Jan '20                        Jan '21

# The Deep Learning Revolution



**Translation**

Google's AI translation system is approaching human-level accuracy

**Reading Comprehension**

ALIBABA AI BEATS HUMANS COMPREHENSION TEST
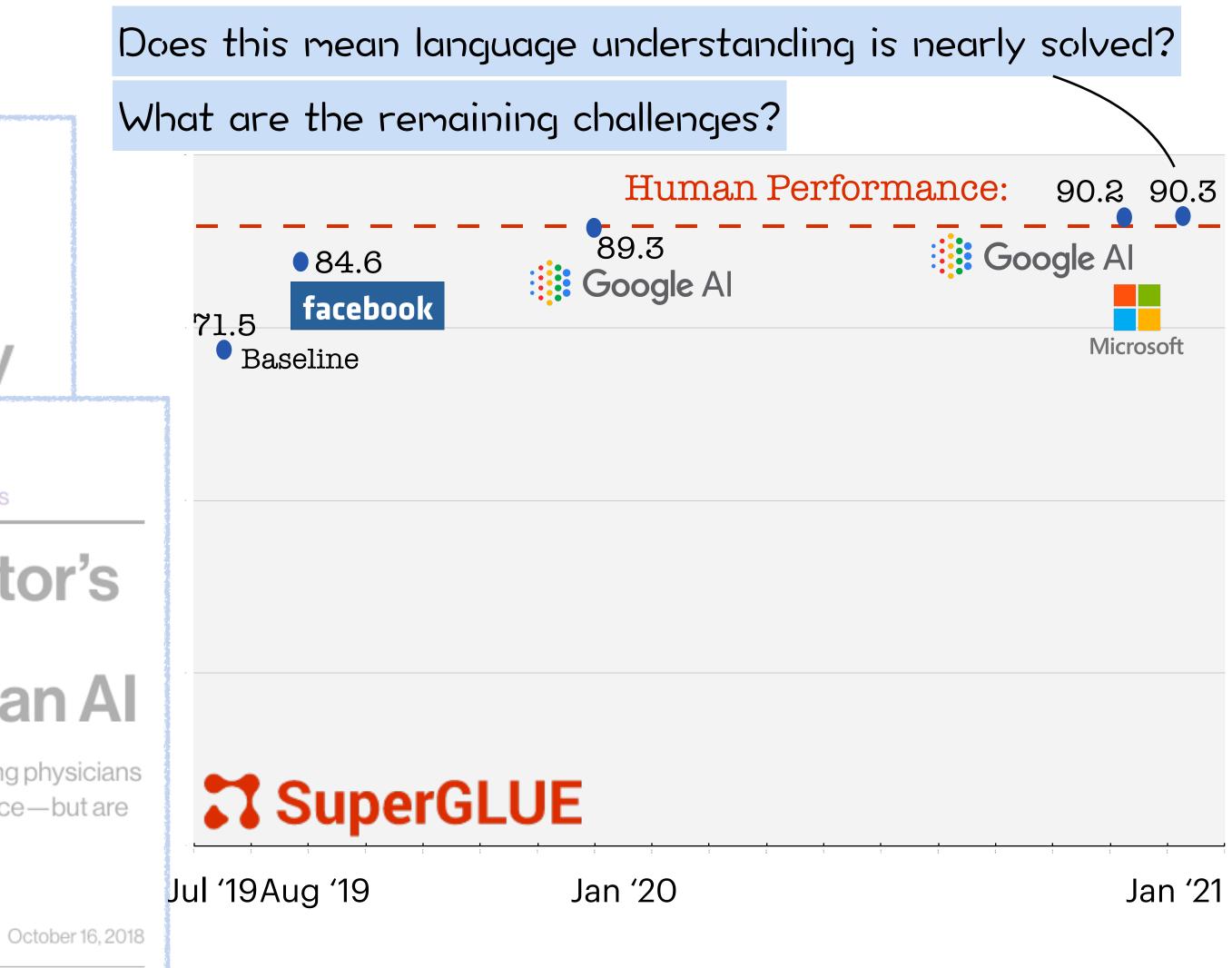
CHRISTINE CHOU | JULY 9, 2019

**Chatbots**

Artificial intelligence / Voice assistants

Your next doctor's appointment might be with an AI

A new wave of chatbots are replacing physicians and providing frontline medical advice—but are they as good as the real thing?

by Will Douglas Heaven          October 16, 2018
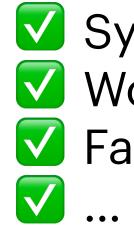
Does this mean language understanding is nearly solved?

What are the remaining challenges?

Human Performance:     90.2   90.3

84.6
facebook

89.3
Google AI

Google AI

Microsoft

71.5
Baseline

SuperGLUE

Jul '19  Aug '19          Jan '20                    Jan '21

# Is Natural Language Understanding Nearly Solved?

**Pre-training**

# Is Natural Language Understanding Nearly Solved?

**Pre-training**



✅ Syntax
✅ Word meanings
✅ Factual Knowledge
✅ ...

# Is Natural Language Understanding Nearly Solved?

**Pre-training**



**Fine-tuning:**

Language Model

✅ Syntax
✅ Word meanings
✅ Factual Knowledge
✅ ...

# Is Natural Language Understanding Nearly Solved?

**Pre-training**



Google AI

WIKIPEDIA
The Free Encyclopedia

→ **Fine-tuning:**

Language Model

↑ ↑ ↑ ↑ ↑

The chocolate cake is amazing

✅ Syntax
✅ Word meanings
✅ Factual Knowledge
✅ ...

# Is Natural Language Understanding Nearly Solved?

**Pre-training**

Google AI

WIKIPEDIA
The Free Encyclopedia

http://www

→ **Fine-tuning:**

5.4%  94.6%

− +

Language Model

↑    ↑    ↑    ↑    ↑
The chocolate cake is amazing

✅ Syntax
✅ Word meanings
✅ Factual Knowledge
✅ ...

# Is Natural Language Understanding Nearly Solved?

**Pre-training**



✅ Syntax
✅ Word meanings
✅ Factual Knowledge
✅ ...

**Fine-tuning:**

5.4%  94.6%

Language Model

The chocolate cake is amazing

✅ Understanding the task
✅ Learning to solve the task

# Is Natural Language Understanding Nearly Solved?

**Pre-training**



✅ Syntax
✅ Word meanings
✅ Factual Knowledge
✅ ...

**Fine-tuning:**

5.4%  94.6%

− +

Language Model

The chocolate cake is amazing

✅ Understanding the task
✅ Learning to solve the task

What are the remaining challenges?

❓ Generalization to unknown situations

# Overfitting to Data-specific Spurious Correlations

*Analyzing the Behavior of Visual Question Answering Models*. Aishwarya Agrawal, Dhruv Batra, and Devi Parikh. EMNLP 2016.

# Overfitting to Data-specific Spurious Correlations

How many zebras?

🤖: 2

*Analyzing the Behavior of Visual Question Answering Models*. Aishwarya Agrawal, Dhruv Batra, and Devi Parikh. EMNLP 2016.

# Overfitting to Data-specific Spurious Correlations

How many zebras?



🤖: 2


How many giraffes? 2


How many zebras? 2


How many dogs? 2

*Analyzing the Behavior of Visual Question Answering Models.* Aishwarya Agrawal, Dhruv Batra, and Devi Parikh. EMNLP 2016.

# Overfitting to Data-specific Spurious Correlations

How many zebras?



🤖: 2

How many giraffes? 2

How many zebras? 2

How many dogs? 2

## ...Solving *datasets* but not underlying *tasks*!

# Humans generalize from few examples

# Humans generalize from few examples

# Humans generalize from few examples

# Humans generalize from few examples

# Humans generalize from few examples

# Commonsense Reasoning
## in Natural Language Processing

# Commonsense Reasoning

# Commonsense Reasoning

**Natural language is...**

# Commonsense Reasoning

**Natural language is...**

**Ambiguous**



Stevie Wonder announces he'll be having kidney surgery during London concert

# Commonsense Reasoning

## Natural language is...

### Ambiguous



Stevie Wonder announces he'll be having kidney surgery during London concert

**Q:** When is the surgery?
**A:** During London concert ❌

# Commonsense Reasoning

## Natural language is...
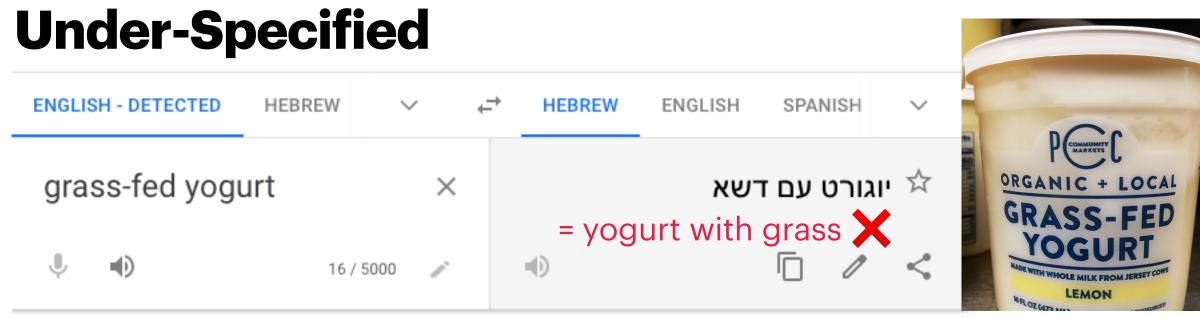
### Ambiguous



Stevie Wonder announces he'll be having kidney surgery during London concert

**Q:** When is the surgery?
**A:** During London concert ❌

🤔 Kidney surgery is performed under general anesthesia
🤔 People are unconscious under general anesthesia
🤔 Performing actions requires being conscious

# Commonsense Reasoning

## Natural language is…

### Ambiguous

Stevie Wonder announces he'll be having kidney surgery during London concert

**Q:** When is the surgery?
**A:** During London concert ❌

🤔 Kidney surgery is performed under general anesthesia
🤔 People are unconscious under general anesthesia
🤔 Performing actions requires being conscious

### Under-Specified

ENGLISH - DETECTED   HEBREW   ⌄   ⇄   HEBREW   ENGLISH   SPANISH   ⌄

grass-fed yogurt   ✕                              ☆ יוגורט עם דשא

🎤  🔊                16 / 5000  ✎      🔊         = yogurt with grass ❌   📋  ✎  ⤴

# Commonsense Reasoning

## Natural language is...

### Ambiguous

Stevie Wonder announces he'll be having kidney surgery during London concert

**Q:** When is the surgery?
**A:** During London concert ❌

🤔 Kidney surgery is performed under general anesthesia
🤔 People are unconscious under general anesthesia
🤔 Performing actions requires being conscious

### Under-Specified

ENGLISH - DETECTED     HEBREW          ⇄     HEBREW    ENGLISH    SPANISH

grass-fed yogurt          ✕                      יוגורט עם דשא   ☆
                                           = yogurt with grass ❌
🎤  🔊              16 / 5000  ✏️        🔊              

ORGANIC + LOCAL
GRASS-FED YOGURT
MADE WITH WHOLE MILK FROM JERSEY COWS
LEMON

🤔 Yogurt is typically made of cow milk
🤔 Cows eat grass

# What is Commonsense?

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

*Introductory Tutorial on Commonsense Reasoning*. Maarten Sap, **Vered Shwartz**, Antoine Bosselut, Dan Roth, and Yejin Choi. ACL 2020.
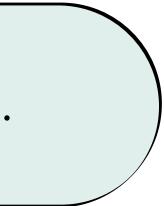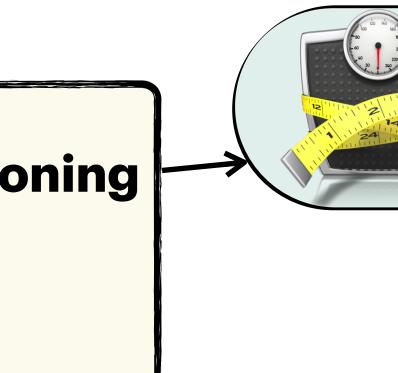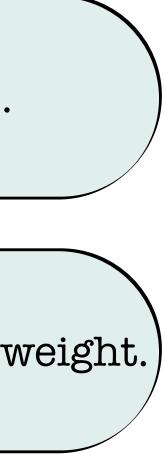
# What is Commonsense?



It's a bad idea to touch a hot stove.

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

https://web.njit.edu/~ronkowit/eliza.html
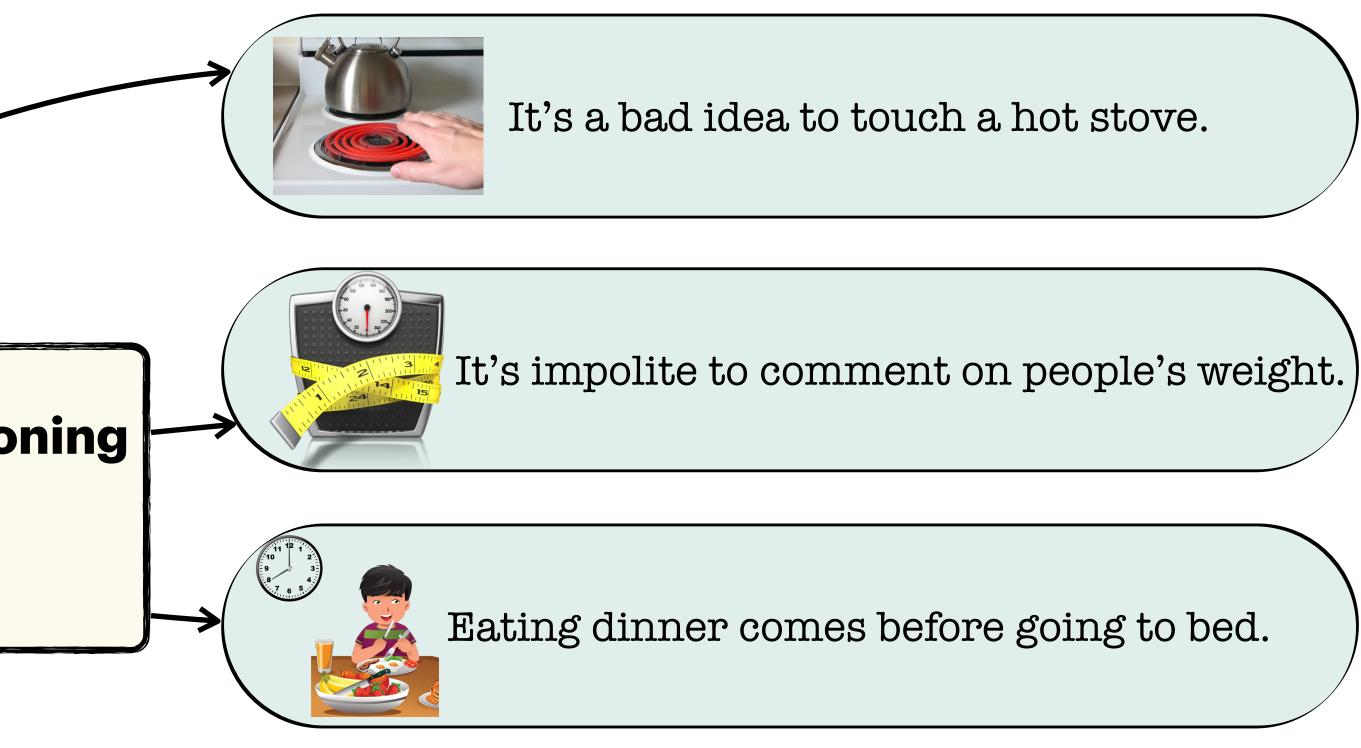
# What is Commonsense?

It's a bad idea to touch a hot stove.

It's impolite to comment on people's weight.

The basic level of **practical knowledge** and **reasoning** concerning **everyday situations** and **events** that are **commonly** shared among **most** people.
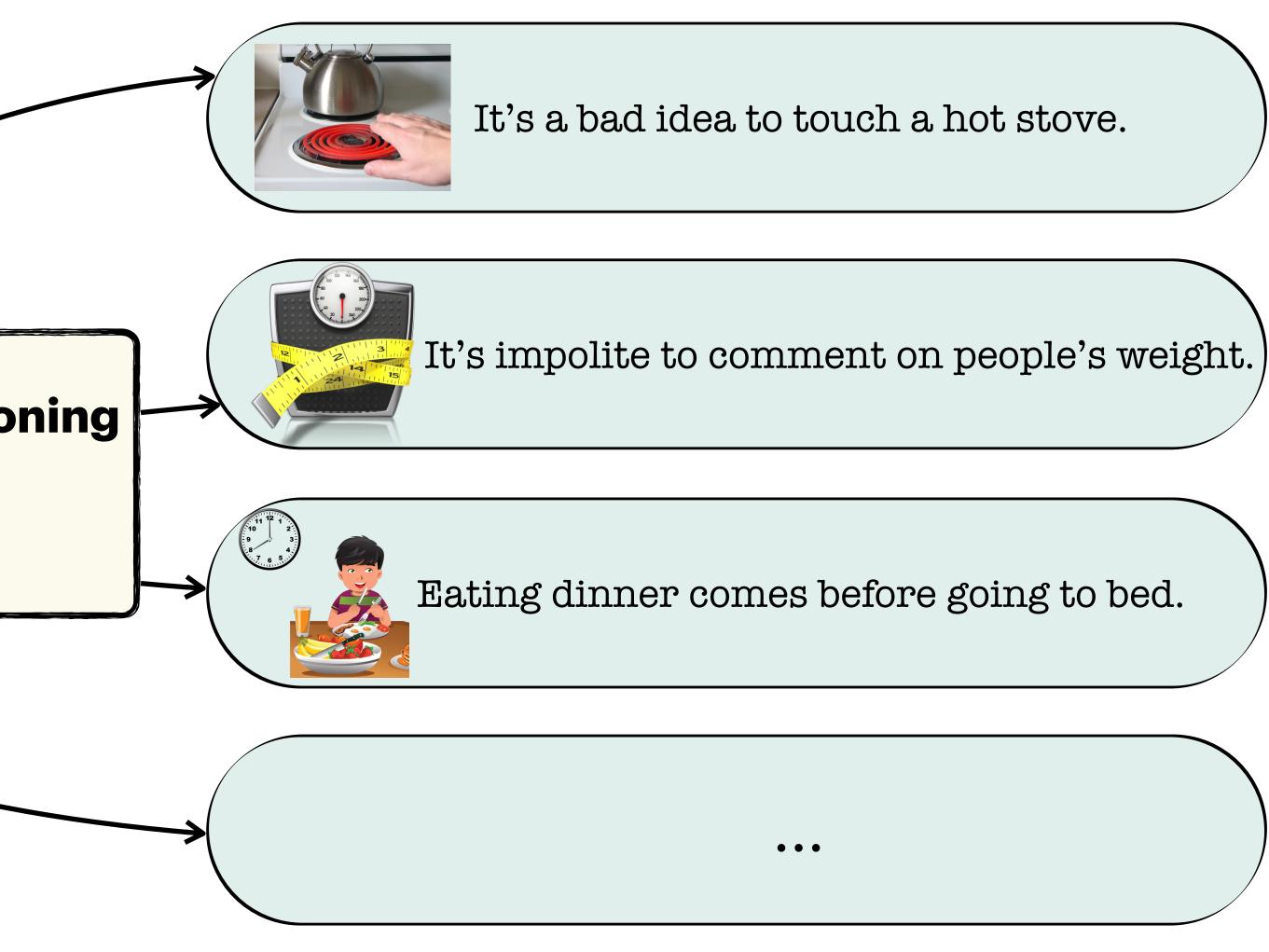
https://web.njit.edu/~ronkowit/eliza.html

# What is Commonsense?

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

It's a bad idea to touch a hot stove.

It's impolite to comment on people's weight.

Eating dinner comes before going to bed.

*Introductory Tutorial on Commonsense Reasoning.* Maarten Sap, **Vered Shwartz**, Antoine Bosselut, Dan Roth, and Yejin Choi. ACL 2020.

https://web.njit.edu/~ronkowit/eliza.html

# What is Commonsense?

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


It's a bad idea to touch a hot stove.


It's impolite to comment on people's weight.


Eating dinner comes before going to bed.

...

https://web.njit.edu/~ronkowit/eliza.html

# Commonsense Timeline



John MacCarthy    Marvin Minsky    Claude Shannon    Ray Solomonoff    Alan Newell

Herbert Simon    Arthur Samuel    Oliver Selfridge    Nathaniel Rochester    Trenchard More
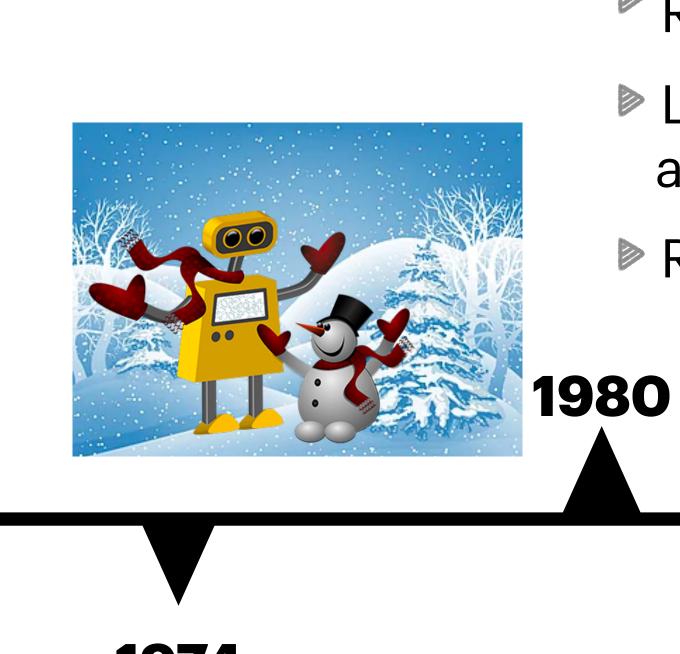
**1956**

# Commonsense Timeline



**1956**

**1974**

File F$_2$ | Directory F$_3$ | Disk F$_4$ | View F

TEXT ⬆

⬇ CALENDAR ⬆

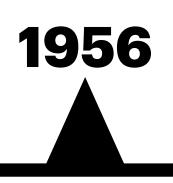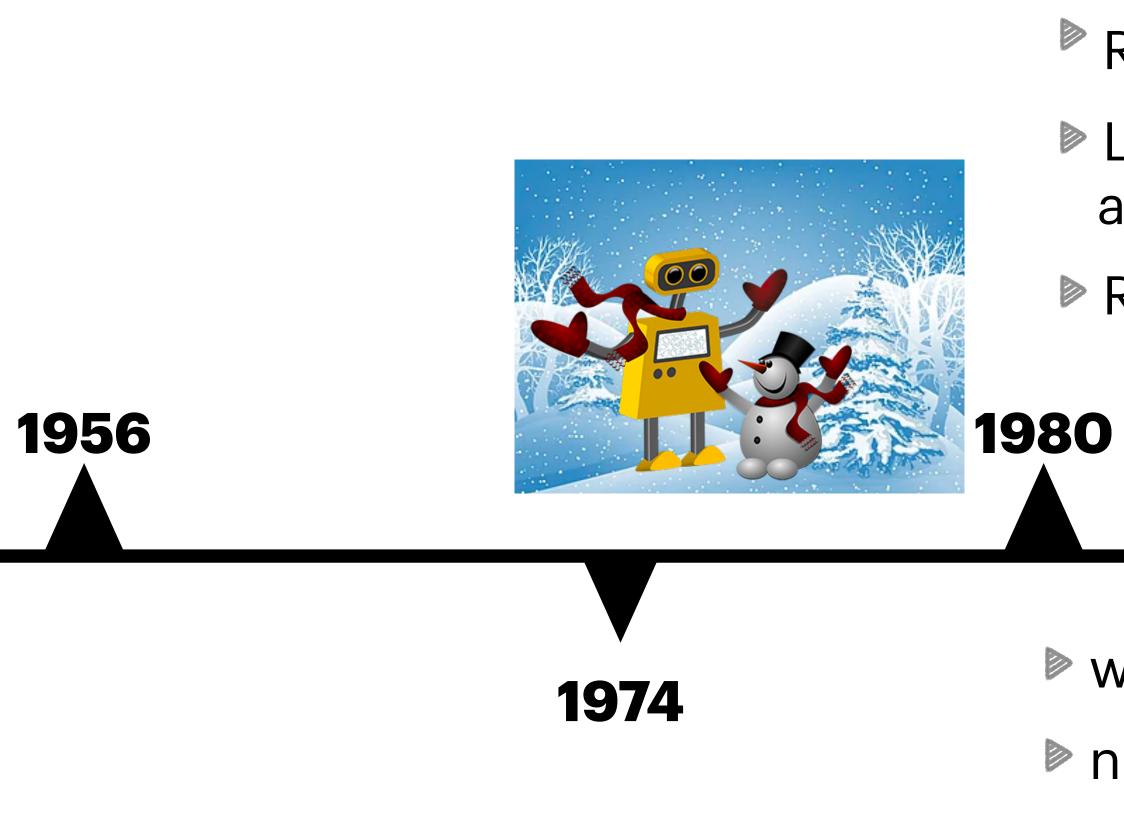File F₂ | Directory F₃ | Disk F₄ | View F

TEXT

CALENDAR

# Commonsense Timeline



▷ Reasoning by search → combinatorial explosion

▷ Lack of commonsense knowledge and reasoning abilities

▷ Rigidity of symbolic reasoning

**1956**
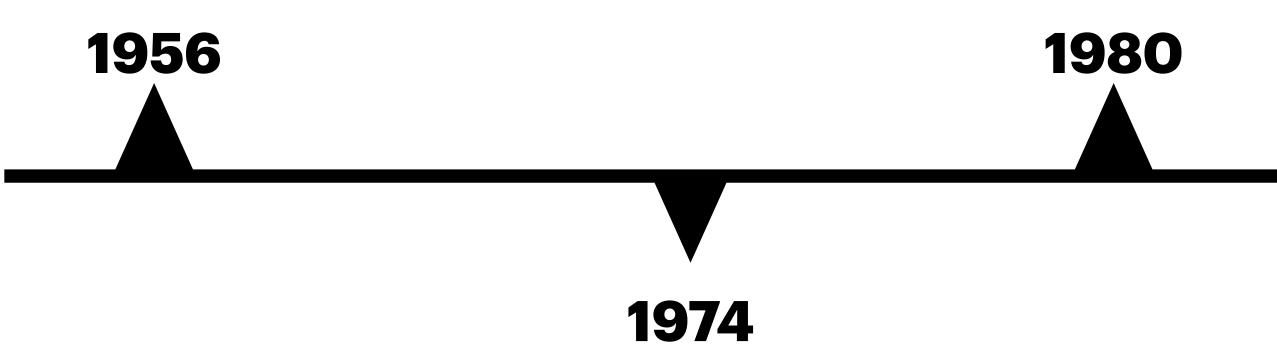
**1980**

**1974**

# Commonsense Timeline



▷ Reasoning by search → combinatorial explosion

▷ Lack of commonsense knowledge and reasoning abilities

▷ Rigidity of symbolic reasoning

**1956**

**1980**

**1974**

▷ weak computing power

▷ not enough data (and no crowdsourcing)

▷ weaker computational models

# Commonsense Timeline

**1956**

**1980**

**1974**

**2011**

▷ Expert systems

▷ Slow progress

# Commonsense Timeline

**1956**

**1974**

**1980**



Big Data



Deep Learning

Inputs

Hidden Layers

Outputs
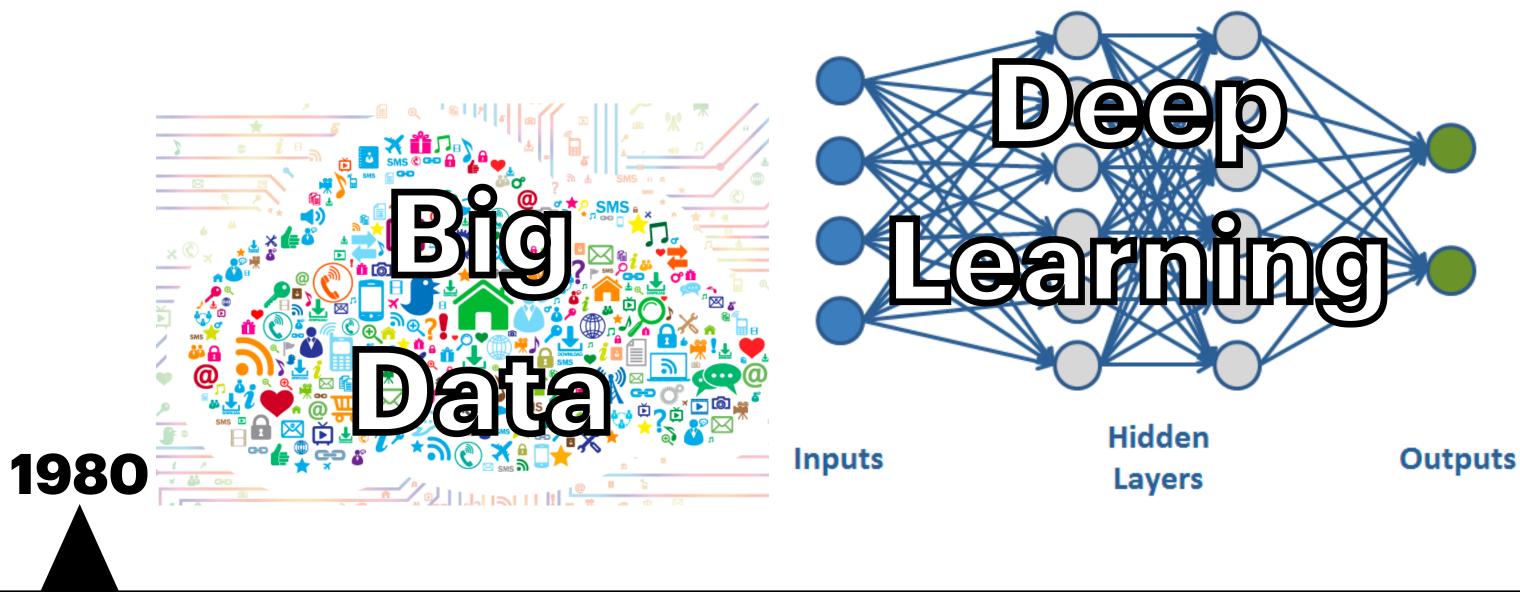
**2011**

# Path to commonsense?

Brute force larger networks with deeper layers?

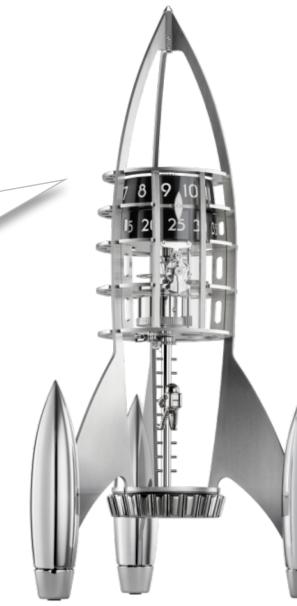# Path to commonsense?

Brute force larger networks with deeper layers?

# Path to commonsense?

Brute force larger networks with deeper layers?

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

# Path to commonsense

Benchmarks

Symbolic Knowledge

Neural Representations

Reasoning engine with commonsense

# Path to commonsense

# 1950: Turing Test


**Alan Turing**

~~Can machines think?~~

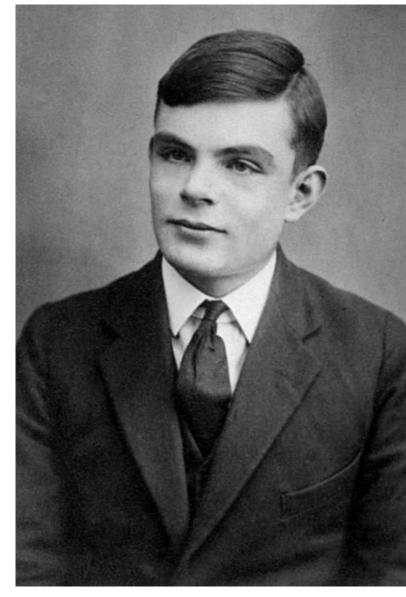Can a human judge distinguish between a human and a machine following a short conversation with each?

# 1950: Turing Test

~~Can machines think?~~

Can a human judge distinguish between a human and a machine following a short conversation with each?
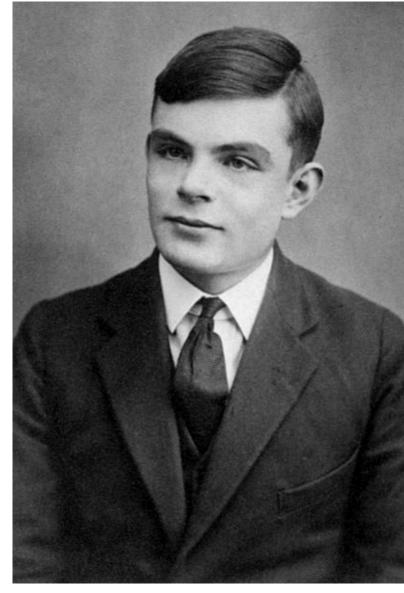


**Alan Turing**

- Loebner Prize (since 1990s)
- Winner of 2014: a bot named "Eugene Goostman", simulating a 13-year-old Ukrainian boy, won
- Recommended reading: https://artistdetective.wordpress.com/, "The most human human"

# Winograd Schema Challenge (WSC)

The winograd schema challenge. Hector Levesque, Ernest Davis, and Leora Morgenstern. AAAI 2012.

# Winograd Schema Challenge (WSC)

The city councilmen refused the demonstrators a permit because *they* **advocated** violence. Who is "*they*"?

(a) The city councilmen
(b) The demonstrators

The winograd schema challenge. Hector Levesque, Ernest Davis, and Leora Morgenstern. AAAI 2012.

# Winograd Schema Challenge (WSC)

The city councilmen refused the demonstrators a permit because *they* **advocated** violence. Who is "*they*"?

(a) The city councilmen
(b) The demonstrators

The winograd schema challenge. Hector Levesque, Ernest Davis, and Leora Morgenstern. AAAI 2012.

# Winograd Schema Challenge (WSC)

The city councilmen refused the demonstrators a permit because *they* **advocated** violence. Who is "*they*"?

(a) The city councilmen
(b) The demonstrators

The city councilmen refused the demonstrators a permit because *they* **feared** violence. Who is "*they*"?
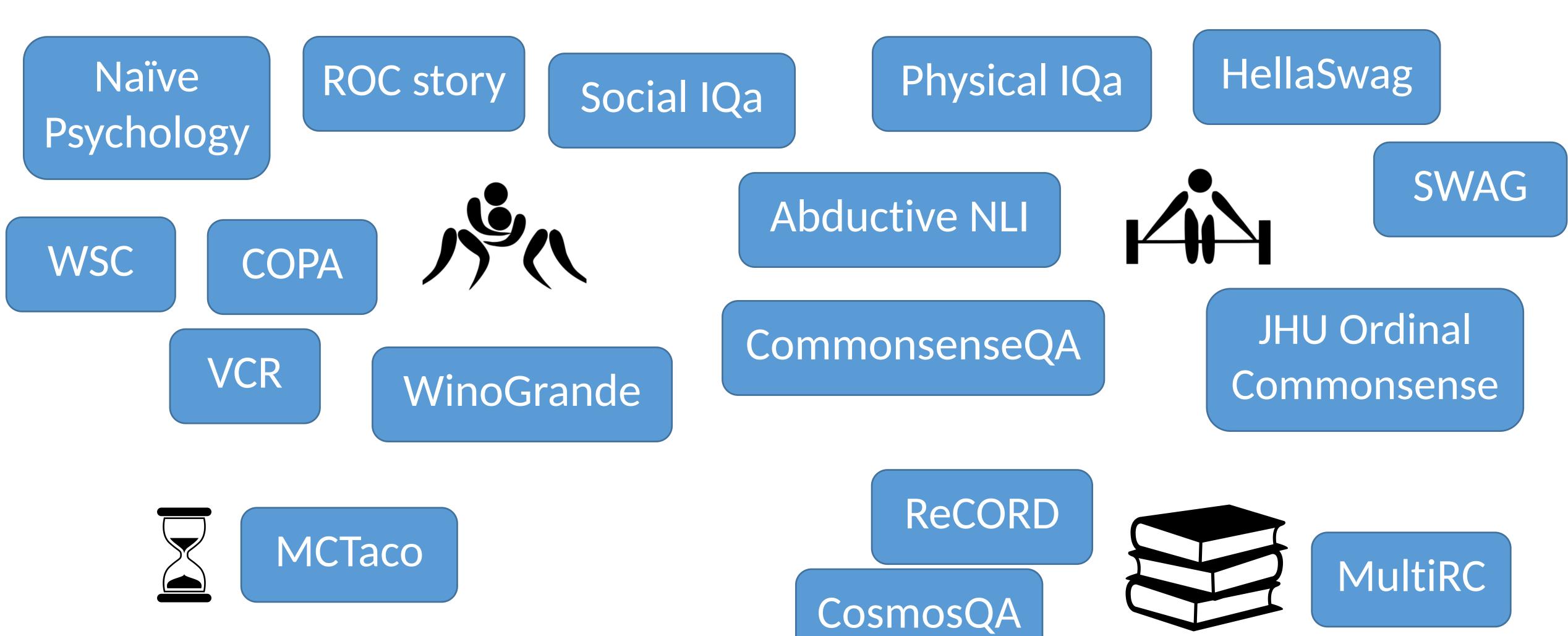
(a) The city councilmen
(b) The demonstrators

The winograd schema challenge. Hector Levesque, Ernest Davis, and Leora Morgenstern. AAAI 2012.

# Winograd Schema Challenge (WSC)

The city councilmen refused the demonstrators a permit because *they* **advocated** violence. Who is "*they*"?
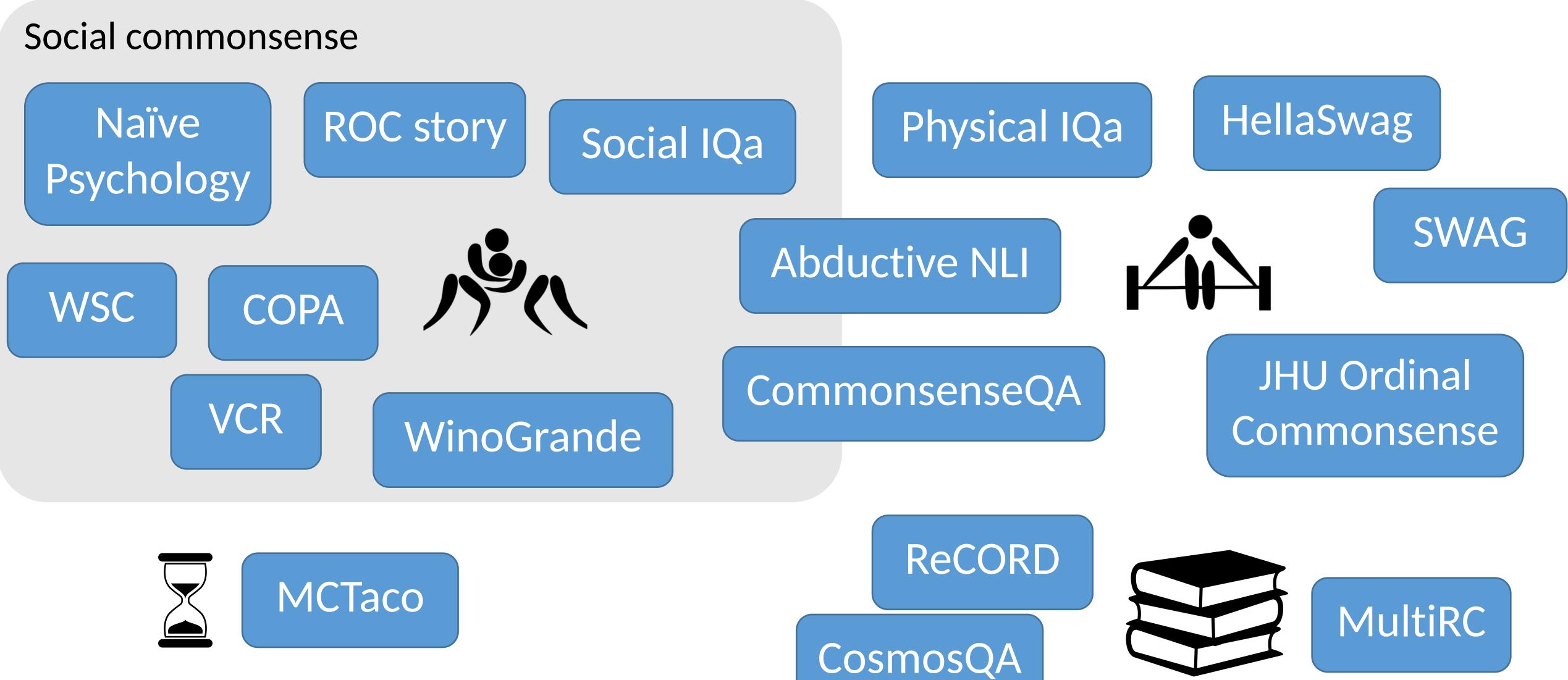
(a) The city councilmen
(b) The demonstrators

The city councilmen refused the demonstrators a permit because *they* **feared** violence. Who is "*they*"?
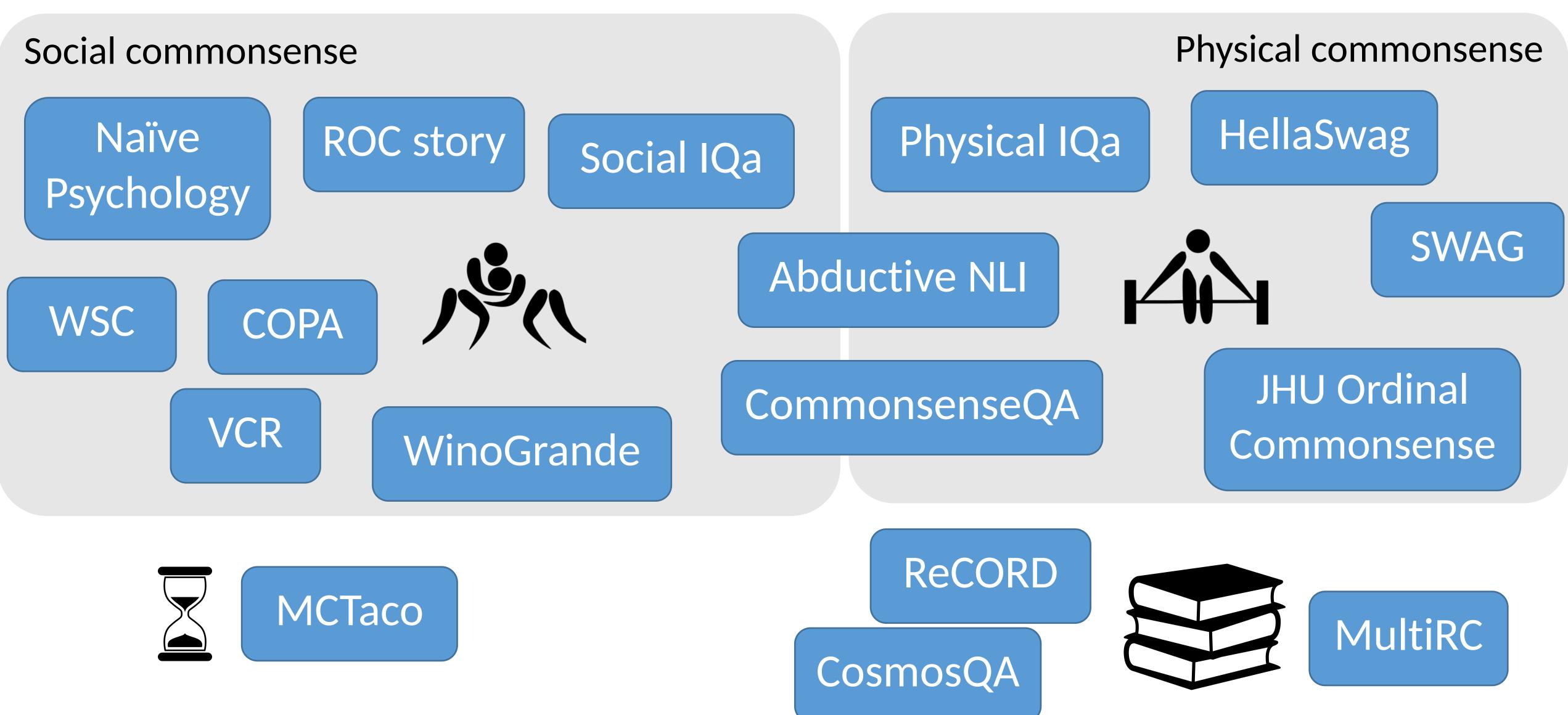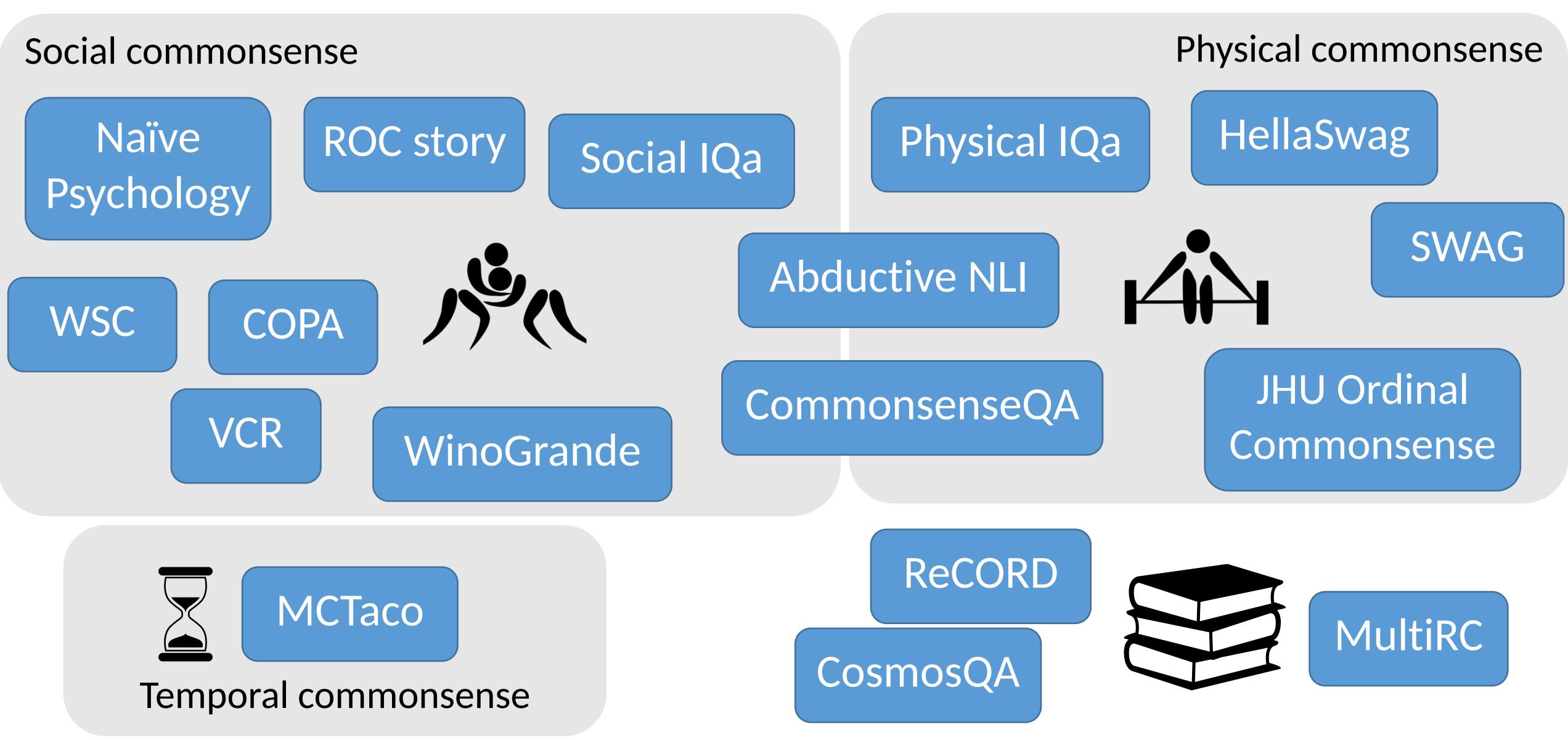
(a) The city councilmen
(b) The demonstrators

The winograd schema challenge. Hector Levesque, Ernest Davis, and Leora Morgenstern. AAAI 2012.

# More benchmarks

Naïve Psychology

ROC story

Social IQa

Physical IQa

HellaSwag

SWAG

WSC

COPA

Abductive NLI

CommonsenseQA

JHU Ordinal Commonsense

VCR

WinoGrande

MCTaco

ReCORD

CosmosQA

MultiRC

Even more benchmarks: https://commonsense.run/

# More benchmarks

Social commonsense

Naïve Psychology

ROC story

Social IQa

Physical IQa

HellaSwag

SWAG

WSC

COPA

Abductive NLI

VCR

WinoGrande

CommonsenseQA

JHU Ordinal Commonsense

MCTaco

ReCORD

CosmosQA

MultiRC

Even more benchmarks: https://commonsense.run/

# More benchmarks

## Social commonsense

Naïve Psychology

ROC story

Social IQa

WSC

COPA

VCR

WinoGrande

## Physical commonsense

Physical IQa

HellaSwag

SWAG

Abductive NLI

CommonsenseQA

JHU Ordinal Commonsense

MCTaco

ReCORD

CosmosQA

MultiRC

Even more benchmarks: https://commonsense.run/

# More benchmarks

**Social commonsense**

Naïve Psychology

ROC story

Social IQa

WSC

COPA

VCR

WinoGrande

**Physical commonsense**

Physical IQa

HellaSwag

SWAG

Abductive NLI

CommonsenseQA

JHU Ordinal Commonsense

MCTaco

Temporal commonsense

ReCORD

CosmosQA

MultiRC

Even more benchmarks: https://commonsense.run/

# More benchmarks

## Social commonsense

Naïve Psychology

ROC story

Social IQa

WSC

COPA

VCR

WinoGrande

## Physical commonsense

Physical IQa

HellaSwag

SWAG

Abductive NLI

CommonsenseQA

JHU Ordinal Commonsense

## Temporal commonsense

MCTaco

## Commonsense reading comprehension

ReCORD

CosmosQA

MultiRC

Even more benchmarks: https://commonsense.run/

# Reasoning about Social Situations

Social IQa

# Reasoning about Social Situations

Alex spilt food all over the floor and it made a huge mess.

What will Alex want to do next?

run around in the mess

mop up the mess

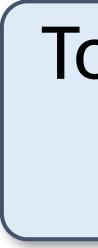# Reasoning about Social Situations

Alex spilt food all over the floor and it made a huge mess.

What will Alex want to do next?

run around in the mess

*less likely*

mop up the mess

*more likely*

**Reasoning about Physical Properties of the World**

To separate egg whites from the yolk using a water bottle, you should

www.youtube.com › watch ▾
Separating Egg Yolks With A Water Bottle - YouTube
EZTV ONLINE is the "How To" channel that combines entertainment with information. We'll show you the ...
Oct. 19, 2015 · Uploaded by eztv online
0:50

Squeeze the water bottle and press it against the yolk. **Release**, which creates suction and lifts the yolk.

Place the water bottle and press it against the yolk. **Keep pushing**, which creates suction and lifts the yolk.

*less likely*          *more likely*

https://leaderboard.allenai.org/physicaliqa/

# COPA: Choice of Plausible Alternatives



The man broke his toe.

What was the cause?

He got a hole in his sock.

*less likely*

He dropped a hammer on his foot.

*more likely*

# RocStories

Karen was assigned a roommate her first year of college. Her roommate asked her to go to a nearby city for a concert. Karen agreed happily. The show was absolutely exhilarating.

Karen hated her roommate.

*less likely*

Karen became good friends with her roommate.

*more likely*

# Discussion: Advantages and Disadvantages of Multiple-Choice Benchmarks

# Reliable Evaluation

# Reliable Evaluation

**Discriminative tasks:**

# Reliable Evaluation

**Discriminative tasks:**

✓ Easy to evaluate

# Reliable Evaluation



**Discriminative tasks:**

✓ Easy to evaluate

✗ Models are right for the wrong reasons

# Reliable Evaluation



**Discriminative tasks:**



**Generative tasks:**

✓ Easy to evaluate

✗ Models are right for the wrong reasons

# Reliable Evaluation

**Discriminative tasks:**

✔️ Easy to evaluate

❌ Models are right for the wrong reasons

**Generative tasks:**

✔️ More nuanced & flexible than pre-defined labels

# Reliable Evaluation

**Discriminative tasks:**

✔️ Easy to evaluate

❌ Models are right for the wrong reasons

**Generative tasks:**

✔️ More nuanced & flexible than pre-defined labels

✔️ More similar to human reasoning process (no "answer choices")

# Reliable Evaluation

**Discriminative tasks:**

✔️ Easy to evaluate

❌ Models are right for the wrong reasons

**Generative tasks:**

✔️ More nuanced & flexible than pre-defined labels

✔️ More similar to human reasoning process (no "answer choices")

✔️ Infinite answer space (no "guessing" of correct answer)

# Reliable Evaluation

**Discriminative tasks:**

✔️ Easy to evaluate

❌ Models are right for the wrong reasons

**Generative tasks:**

✔️ More nuanced & flexible than pre-defined labels

✔️ More similar to human reasoning process (no "answer choices")

✔️ Infinite answer space (no "guessing" of correct answer)

❌ No reliable automatic evaluation metric

# CommonGen



**Concept-Set:** a collection of objects/actions.

dog | frisbee| catch | throw

*Generative Commonsense Reasoning*

**Expected Output:** everyday scenarios covering all given concepts.

- A dog leaps to catch a thrown frisbee.      **[Humans]**
- The dog catches the frisbee when the boy throws it.
- A man throws away his dog 's favorite frisbee expecting him to catch it in the air.

GPT2: A dog throws a frisbee at a football player.      **[Machines]**
UniLM: Two dogs are throwing frisbees at each other .
BART: A dog throws a frisbee and a dog catches it.
T5: dog catches a frisbee and throws it to a dog

**https://inklab.usc.edu/CommonGen/**

# Path to commonsense

Benchmarks

**Symbolic Knowledge**

Neural Representations

Reasoning engine with commonsense

# Grandma's glasses



Tom's grandma was reading a new book, when she dropped her glasses.

She couldn't pick them up, so she called Tom for help.

Tom rushed to help her look for them, they heard a loud crack.

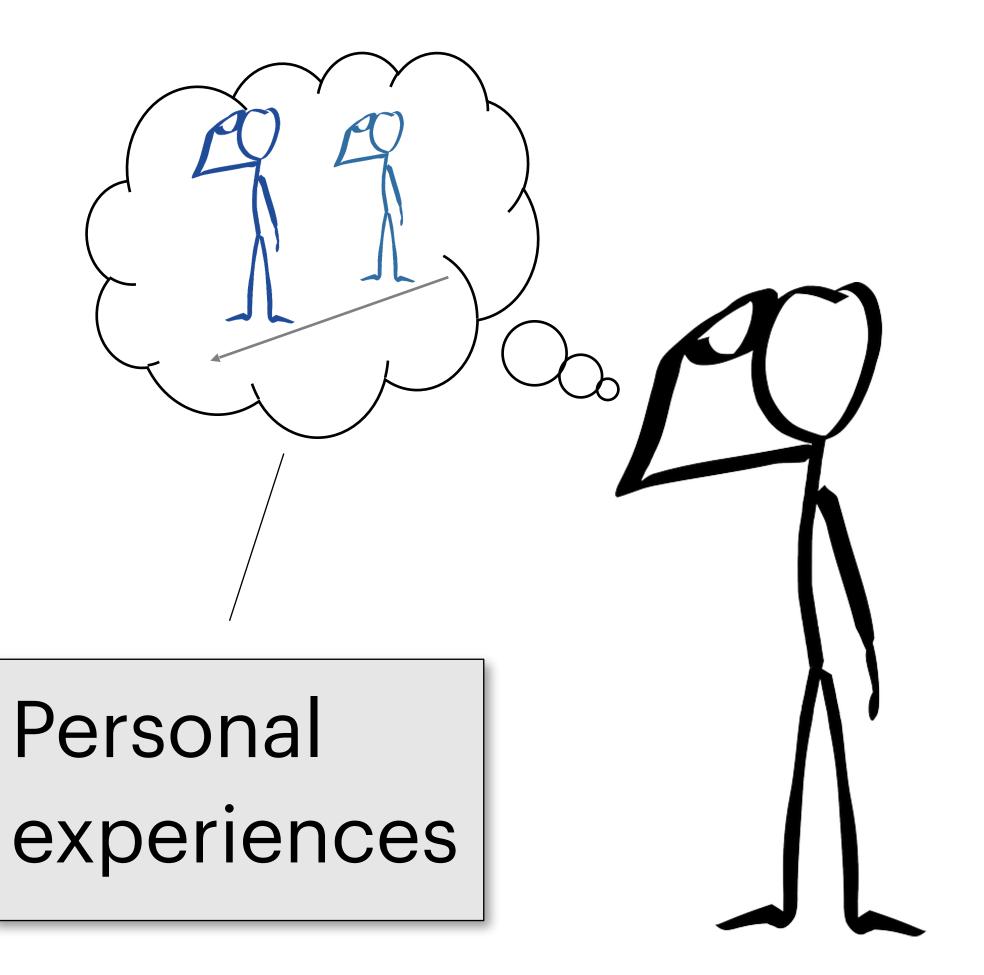They realized that Tom broke her glasses by stepping on them.

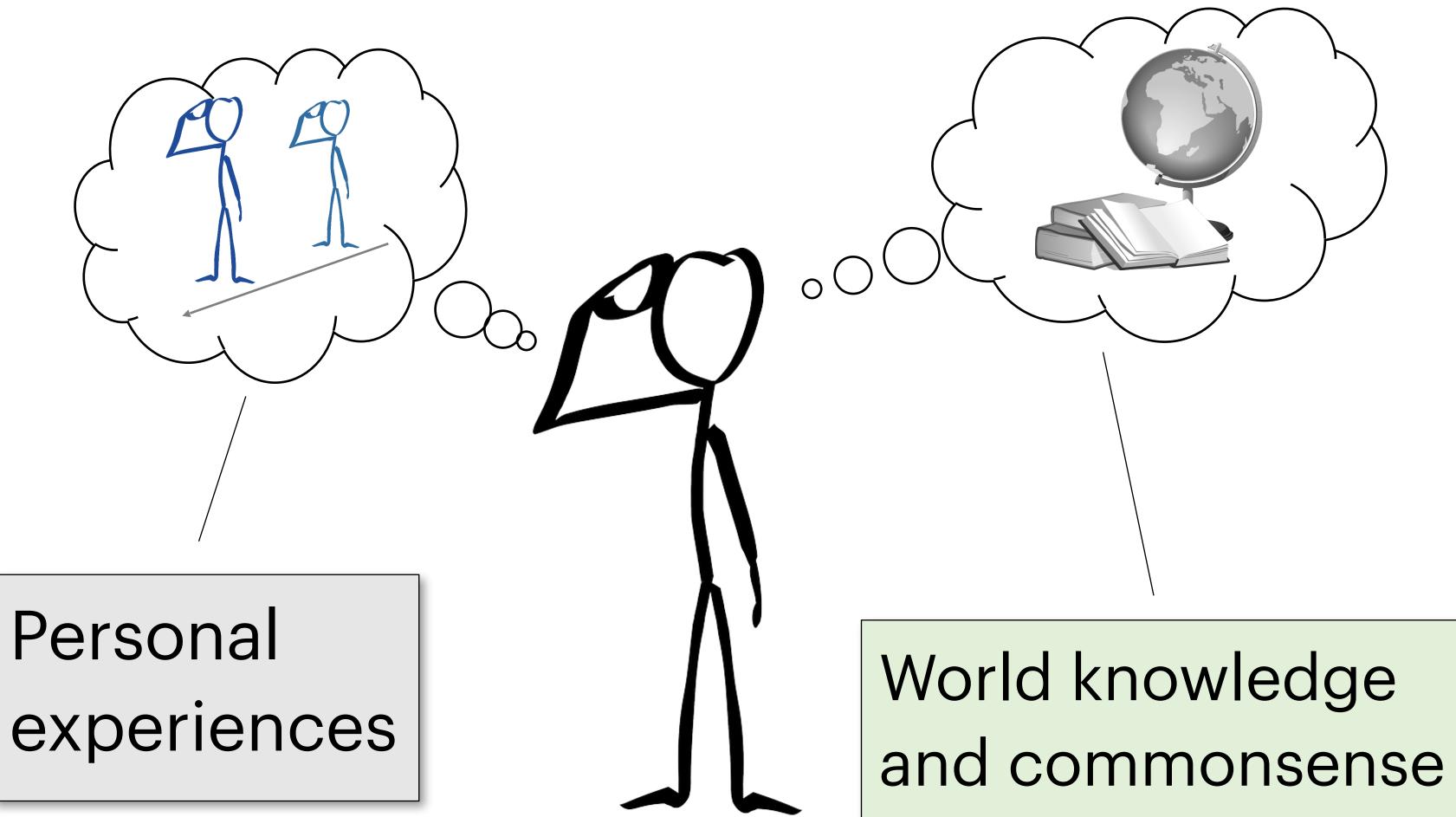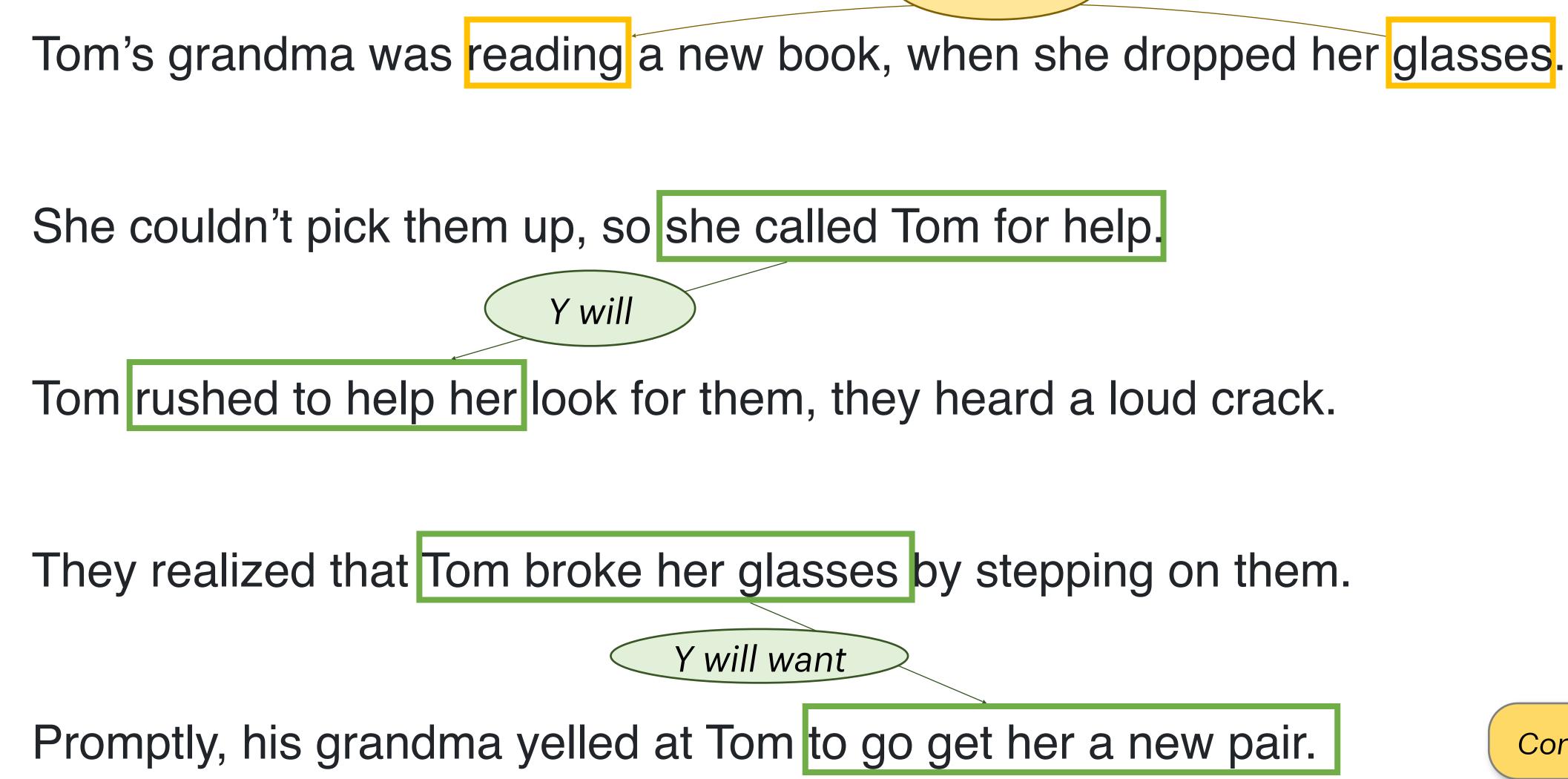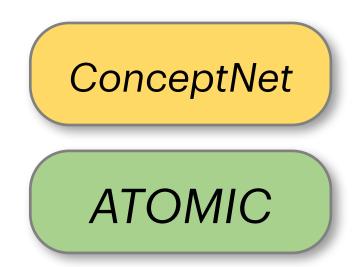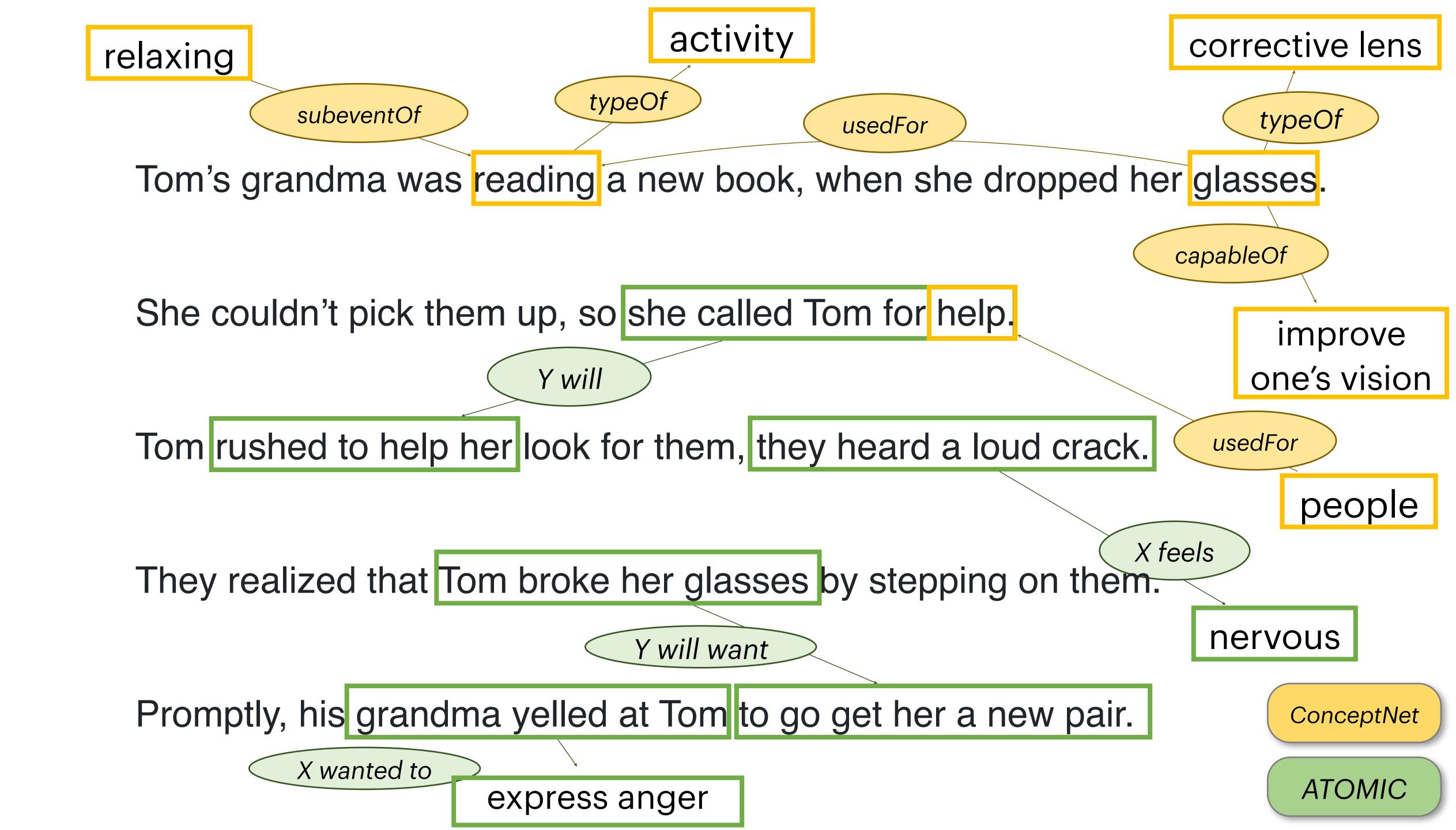Promptly, his grandma yelled at Tom to go get her a new pair.

Humans reason about the world with mental models [Graesser, 1994]

Humans reason about the world with
mental models [Graesser, 1994]



Personal experiences

# Humans reason about the world with mental models [Graesser, 1994]



Personal experiences

World knowledge and commonsense

# Humans reason about the world with mental models [Graesser, 1994]



Personal experiences

World knowledge and commonsense

Commonsense resources aim to be a bank of knowledge for machines to be able to reason about the world in tasks

Tom's grandma was reading a new book, when she dropped her glasses.

She couldn't pick them up, so she called Tom for help.

Tom rushed to help her look for them, they heard a loud crack.

They realized that Tom broke her glasses by stepping on them.

Promptly, his grandma yelled at Tom to go get her a new pair.

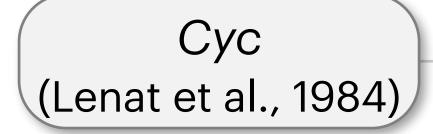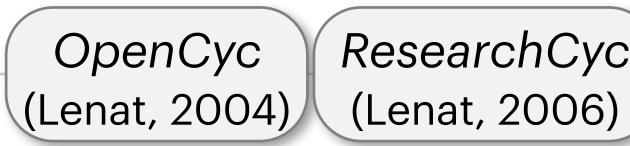Tom's grandma was reading a new book, when she dropped her glasses.

*usedFor*

She couldn't pick them up, so she called Tom for help.

*Y will*

Tom rushed to help her look for them, they heard a loud crack.

They realized that Tom broke her glasses by stepping on them.

*Y will want*

Promptly, his grandma yelled at Tom to go get her a new pair.

*ConceptNet*

*ATOMIC*

relaxing

activity

corrective lens

*subeventOf*

*typeOf*

*usedFor*

*typeOf*

Tom's grandma was reading a new book, when she dropped her glasses.

*capableOf*

She couldn't pick them up, so she called Tom for help.

improve one's vision

*Y will*

Tom rushed to help her look for them, they heard a loud crack.

*usedFor*

people

*X feels*

They realized that Tom broke her glasses by stepping on them.

nervous

*Y will want*

Promptly, his grandma yelled at Tom to go get her a new pair.

*X wanted to*

express anger

*ConceptNet*

*ATOMIC*

# Overview of existing resources

Represented in **symbolic logic**
(e.g., LISP-style logic)

```
(#$implies
 (#$and
  (#$isa ?OBJ ?SUBSET)
  (#$genls ?SUBSET ?SUPERSET))
(#$isa ?OBJ ?SUPERSET))
```

*Cyc*
(Lenat et al., 1984)

*OpenCyc*
(Lenat, 2004)

*ResearchCyc*
(Lenat, 2006)

*OpenCyc 4.0*
(Lenat, 2012)

*today*

# Overview of existing resources



*Open Mind Common Sense*
(Minsky, Singh & Havasi, 1999)

*ConceptNet*
(Liu & Singh, 2004)

*ConceptNet 5.5*
(Speer et al., 2017)

*Cyc*
(Lenat et al., 1984)

*OpenCyc*
(Lenat, 2004)
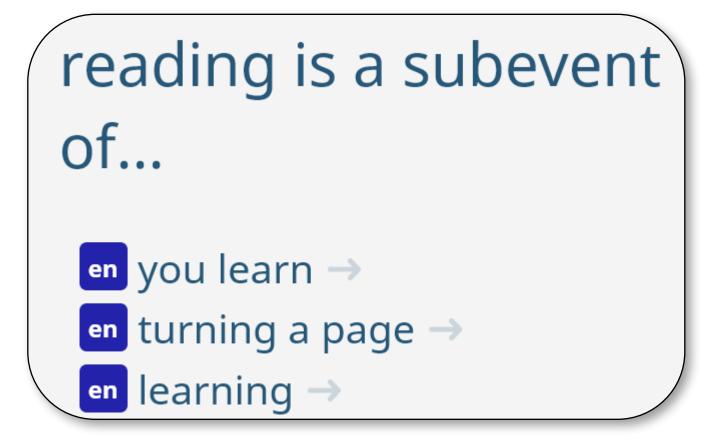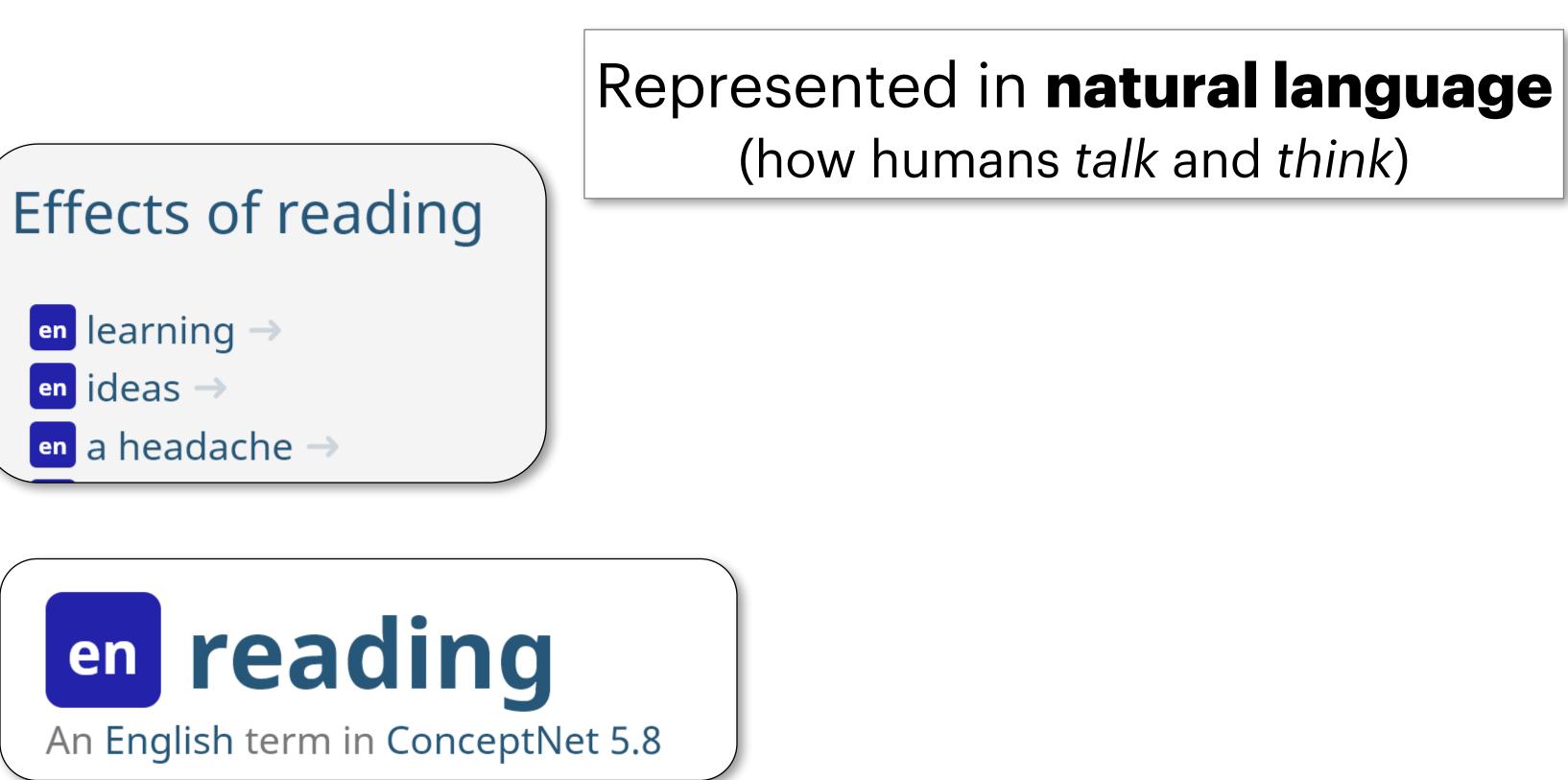
*ResearchCyc*
(Lenat, 2006)

*OpenCyc 4.0*
(Lenat, 2012)

*today*

Represented in **natural language**
(how humans *talk* and *think*)

en reading
An English term in ConceptNet 5.8

reading is a subevent of...

en you learn →

en turning a page →

en learning →

en **reading**

An English term in ConceptNet 5.8

## Related terms

en book →

en books →

en book →

## reading is a subevent of...

en you learn →

en turning a page →

en learning →

en **reading**

An English term in ConceptNet 5.8

Represented in **natural language**
(how humans *talk* and *think*)

## Related terms

- **en** book →
- **en** books →
- **en** book →

## Effects of reading

- **en** learning →
- **en** ideas →
- **en** a headache →

## reading is a subevent of...

- **en** you learn →
- **en** turning a page →
- **en** learning →

**en** **reading**

An English term in ConceptNet 5.8

Represented in **natural language**
(how humans *talk* and *think*)

Related terms

en book →
en books →
en book →

Effects of reading

en learning →
en ideas →
en a headache →

Represented in **natural language**
(how humans *talk* and *think*)

reading is a type of...

en an activity →
en a good way to learn →
en one way of learning →
en one way to learn →

reading is a subevent of...

en you learn →
en turning a page →
en learning →

en reading
An English term in ConceptNet 5.8

## Related terms

- en book →
- en books →
- en book →

## Effects of reading

- en learning →
- en ideas →
- en a headache →

## Represented in **natural language**
(how humans *talk* and *think*)

## reading is a type of...

- en an activity →
- en a good way to learn →
- en one way of learning →
- en one way to learn →

## reading is a subevent of...

- en you learn →
- en turning a page →
- en learning →

## en reading

An English term in ConceptNet 5.8

## Types of reading

- en browse (n, communication) →
- en bumf (n, communication) →
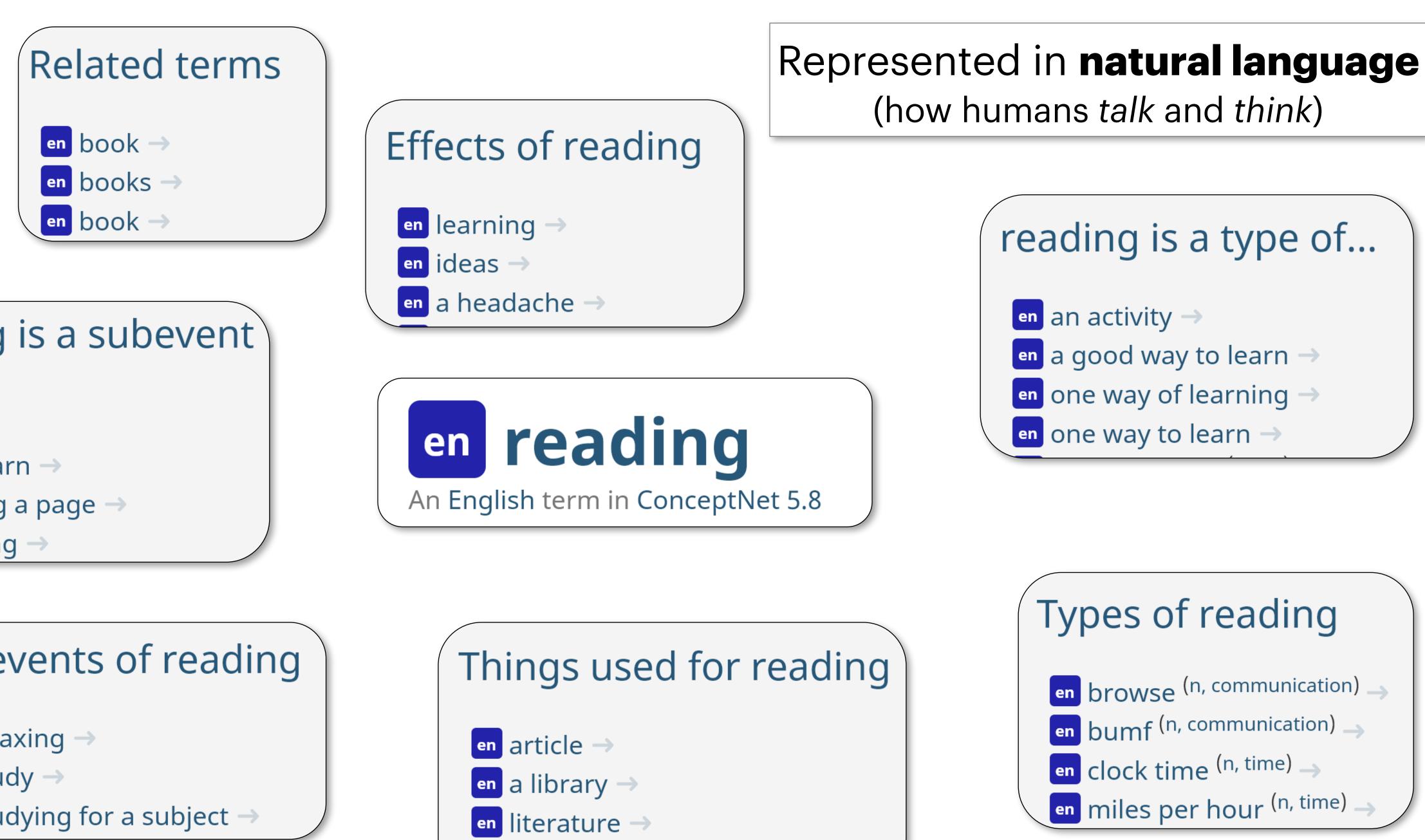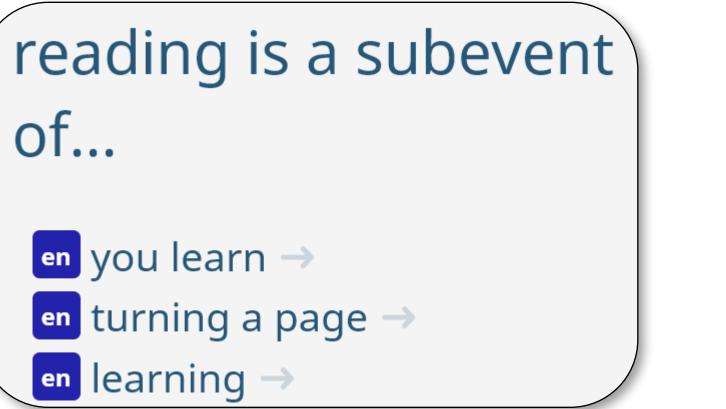- en clock time (n, time) →
- en miles per hour (n, time) →

## Related terms

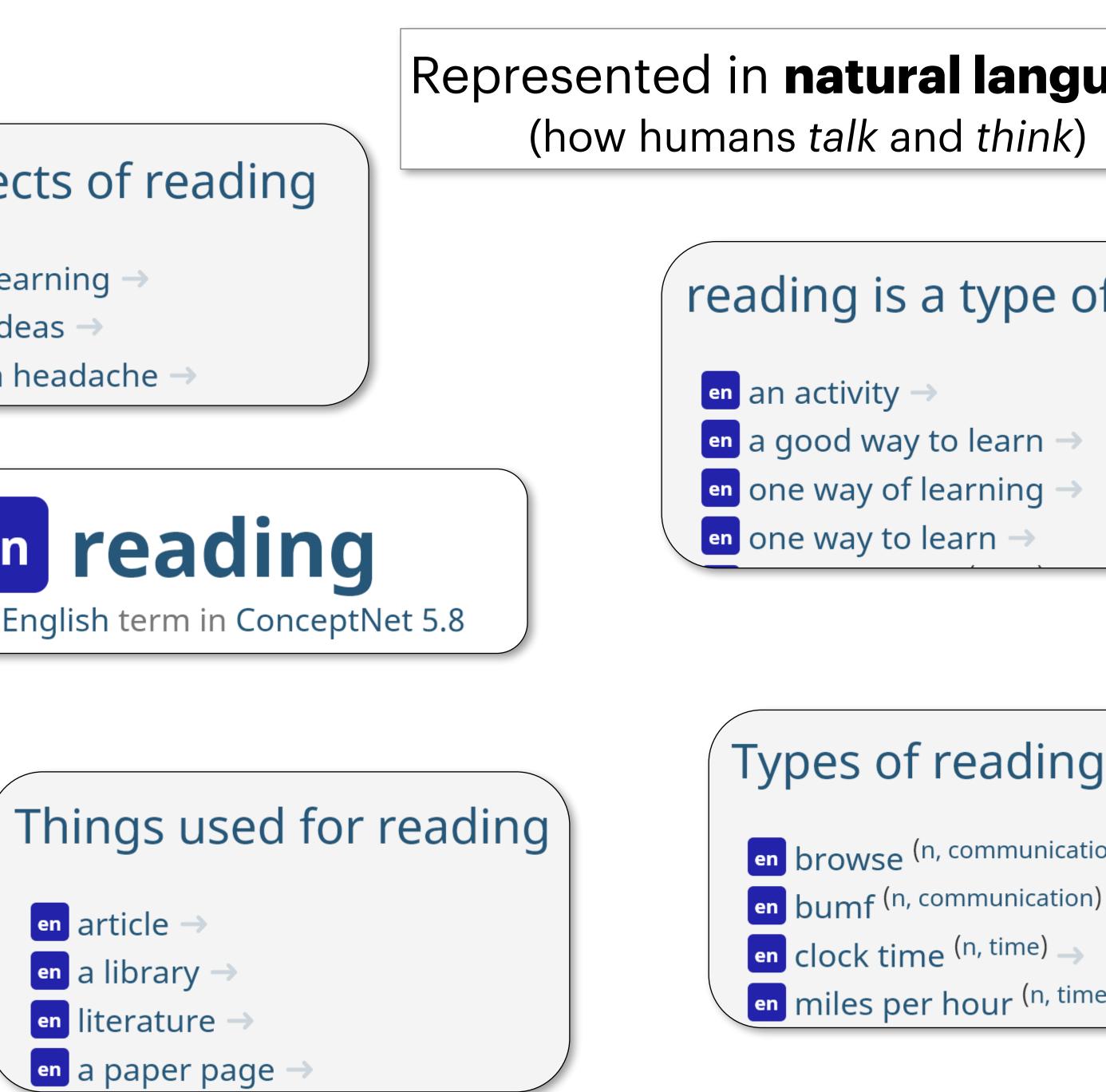**en** book →
**en** books →
**en** book →

## Effects of reading

**en** learning →
**en** ideas →
**en** a headache →

## Represented in **natural language**
(how humans *talk* and *think*)

## reading is a type of...

**en** an activity →
**en** a good way to learn →
**en** one way of learning →
**en** one way to learn →

## reading is a subevent of...

**en** you learn →
**en** turning a page →
**en** learning →

## en reading

An English term in ConceptNet 5.8

## Things used for reading

**en** article →
**en** a library →
**en** literature →
**en** a paper page →

## Types of reading

**en** browse (n, communication) →
**en** bumf (n, communication) →
**en** clock time (n, time) →
**en** miles per hour (n, time) →

Represented in **natural language**
(how humans *talk* and *think*)

Related terms

en book →
en books →
en book →

Effects of reading

en learning →
en ideas →
en a headache →

reading is a type of...

en an activity →
en a good way to learn →
en one way of learning →
en one way to learn →

reading is a subevent of...

en you learn →
en turning a page →
en learning →

en reading
An English term in ConceptNet 5.8

Subevents of reading

en relaxing →
en study →
en studying for a subject →

Things used for reading

en article →
en a library →
en literature →
en a paper page →

Types of reading

en browse (n, communication) →
en bumf (n, communication) →
en clock time (n, time) →
en miles per hour (n, time) →

# Overview of existing resources



NELL
(Carlson et al., 2010)

NELL
(Mitchell et al., 2015)

Open Mind Common Sense
(Minsky, Singh & Havasi, 1999)

ConceptNet
(Liu & Singh, 2004)

ConceptNet 5.5
(Speer et al., 2017)

Cyc
(Lenat et al., 1984)

OpenCyc
(Lenat, 2004)

ResearchCyc
(Lenat, 2006)

OpenCyc 4.0
(Lenat, 2012)

today

# Overview of existing resources

# Overview of existing resources

Represented in **natural language**
(how humans *talk* and *think*)

**ATOMIC:** 880,000 triples for AI systems to reason
about *causes* and *effects* of everyday situations

# Decisions when building a new resource

# Decisions when building a new resource

**1. Representation**  Tradeoff between **expressivity** and **ease of collection**

# Decisions when building a new resource

**1. Representation**  Tradeoff between **expressivity** and **ease of collection**

**2. Knowledge Type**

# Decisions when building a new resource

**1. Representation**  Tradeoff between **expressivity** and **ease of collection**

**2. Knowledge Type**

**3. Acquisition Method**

# Discussion:
# Tradeoffs between collecting knowledge from people and extracting from text

# 3. Acquisition Method

**1️⃣ from people**

**2️⃣ from text**

❌ Expensive, takes a long time 💲💲💲

# 3. Acquisition Method

1️⃣ **from people**

2️⃣ **from text**

✘ Expensive, takes a long time 💲💲

*Reporting bias and knowledge acquisition.* Jonathan Gordon and Benjamin Van Durme. AKBC 2013.

# 3. Acquisition Method

**1 from people**

**2 from text**

✖ Expensive, takes a long time $$$

✖ Reporting bias

*Reporting bias and knowledge acquisition*. Jonathan Gordon and Benjamin Van Durme. AKBC 2013.

# 3. Acquisition Method

**1 from people**

**2 from text**

✖ Expensive, takes a long time $$$

✖ Reporting bias

murdered + killed

breathed + exhaled + inhaled

*Reporting bias and knowledge acquisition.* Jonathan Gordon and Benjamin Van Durme. AKBC 2013.

# 3. Acquisition Method

**1 from people**

✘ Expensive, takes a long time $$$

**2 from text**

✘ Reporting bias

✘ What is NOT true

*Reporting bias and knowledge acquisition*. Jonathan Gordon and Benjamin Van Durme. AKBC 2013.

# Path to commonsense

Benchmarks → Symbolic Knowledge → **Neural Representations** → Reasoning engine with commonsense

# ✔ Knowledge in Pre-trained LMs

*A Primer in BERTology: What we know about how BERT works*. Anna Rogers, Olga Kovaleva, and Anna Rumshisky. TACL 2020.
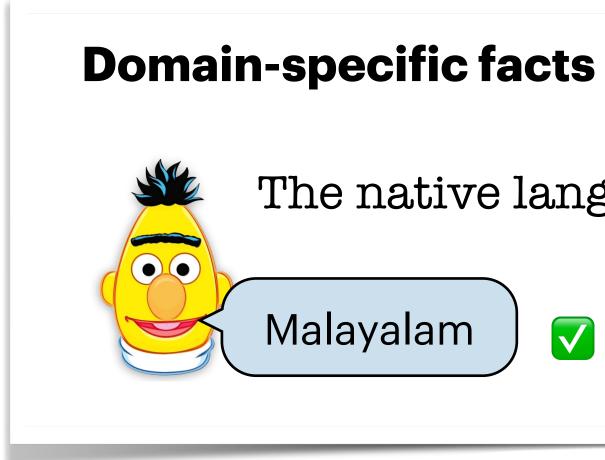
# ✔️ Knowledge in Pre-trained LMs

✅ Syntax:
- Encode information about parts of speech, syntactic chunks and roles
- Syntax trees can be recovered from the representation
- Subject-verb agreement (e.g. tense, plurality)

*A Primer in BERTology: What we know about how BERT works*. Anna Rogers, Olga Kovaleva, and Anna Rumshisky. TACL 2020.

# ✔ **Knowledge in Pre-trained LMs**

✅ Syntax:
- Encode information about parts of speech, syntactic chunks and roles
- Syntax trees can be recovered from the representation
- Subject-verb agreement (e.g. tense, plurality)

✅ Semantics:
- Semantic roles
- Entity types

*A Primer in BERTology: What we know about how BERT works*. Anna Rogers, Olga Kovaleva, and Anna Rumshisky. TACL 2020.

# ✔ Knowledge in Pre-trained LMs



✅ Syntax:
- Encode information about parts of speech, syntactic chunks and roles
- Syntax trees can be recovered from the representation
- Subject-verb agreement (e.g. tense, plurality)

✅ Semantics:
- Semantic roles
- Entity types

✅ Factual knowledge

**Domain-specific facts**   Most people don't know

The native language of Mammootty is [MASK].

Malayalam ✅

*A Primer in BERTology: What we know about how BERT works.* Anna Rogers, Olga Kovaleva, and Anna Rumshisky. TACL 2020.

# ❌ Knowledge in Pre-trained LMs

How can we know what language models know? Zhengbao Jiang, Frank F. Xu, Jun Araki, and Graham Neubig. TACL 2020

Negated and Misprimed Probes for Pretrained Language Models: Birds Can Talk, But Cannot Fly. Nora Kassner and Hinrich Schütze. ACL 2020

What BERT is not: Lessons from a new suite of psycholinguistic diagnostics for language models. Allyson Ettinger. TACL 2020
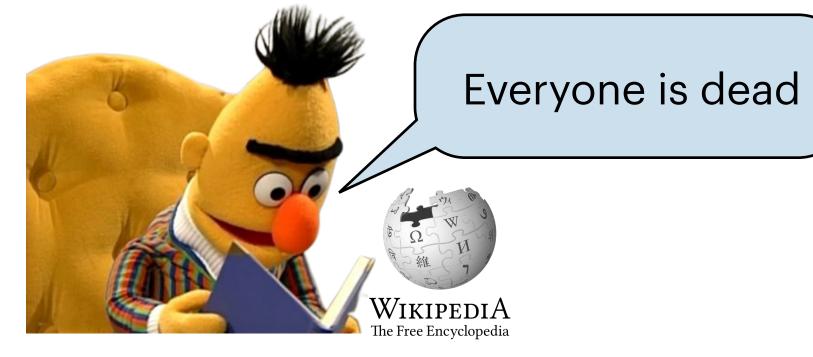
# ❌ Knowledge in Pre-trained LMs

❌ Confuse semantically-similar mutually-exclusive terms



DirectX is developed by [MASK].

| 1 | Intel | -1.06 |
| 2 | Microsoft | -2.21 |
| 3 | IBM | -2.76 |
| 4 | Google | -3.40 |
| 5 | Nokia | -3.58 |

(Jiang et al., 2020)

How can we know what language models know? Zhengbao Jiang, Frank F. Xu, Jun Araki, and Graham Neubig. TACL 2020

Negated and Misprimed Probes for Pretrained Language Models: Birds Can Talk, But Cannot Fly. Nora Kassner and Hinrich Schütze. ACL 2020

What BERT is not: Lessons from a new suite of psycholinguistic diagnostics for language models. Allyson Ettinger. TACL 2020

# ❌ Knowledge in Pre-trained LMs

Birds [MASK] fly.

Can / can't

**❌** Confuse semantically-similar mutually-exclusive terms

**❌** Are really bad with negation

(**Kassner et al. 2020; Ettinger, 2020**)

How can we know what language models know? Zhengbao Jiang, Frank F. Xu, Jun Araki, and Graham Neubig. TACL 2020

Negated and Misprimed Probes for Pretrained Language Models: Birds Can Talk, But Cannot Fly. Nora Kassner and Hinrich Schütze. ACL 2020

What BERT is not: Lessons from a new suite of psycholinguistic diagnostics for language models. Allyson Ettinger. TACL 2020

# ❌ Knowledge in Pre-trained LMs

❌ Confuse semantically-similar mutually-exclusive terms

❌ Are really bad with negation

❌ Lack perceptual knowledge (people don't talk about it)

How can we know what language models know? Zhengbao Jiang, Frank F. Xu, Jun Araki, and Graham Neubig. TACL 2020

Negated and Misprimed Probes for Pretrained Language Models: Birds Can Talk, But Cannot Fly. Nora Kassner and Hinrich Schütze. ACL 2020

What BERT is not: Lessons from a new suite of psycholinguistic diagnostics for language models. Allyson Ettinger. TACL 2020

# ❌ Knowledge in Pre-trained LMs

❌ Confuse semantically-similar mutually-exclusive terms

❌ Are really bad with negation

❌ Lack perceptual knowledge (people don't talk about it)

❌ Also suffer from reporting bias!

How can we know what language models know? Zhengbao Jiang, Frank F. Xu, Jun Araki, and Graham Neubig. TACL 2020

Negated and Misprimed Probes for Pretrained Language Models: Birds Can Talk, But Cannot Fly. Nora Kassner and Hinrich Schütze. ACL 2020

What BERT is not: Lessons from a new suite of psycholinguistic diagnostics for language models. Allyson Ettinger. TACL 2020

# Path to commonsense

# Winograd Schema Challenge (WSC)

The city councilmen refused the demonstrators a permit because *they* **advocated** violence. Who is "*they*"?

(a)The city councilmen
(b)The demonstrators

The city councilmen refused the demonstrators a permit because *they* **feared** violence. Who is "*they*"?

(a)The city councilmen
(b)The demonstrators

The winograd schema challenge. Hector Levesque, Ernest Davis, and Leora Morgenstern. AAAI 2012.

# Supervised Approach

[CLS] The city councilmen refused the demonstrators a permit because [SEP] the city councilmen advocated violence.

[CLS] The city councilmen refused the demonstrators a permit because [SEP] the demonstrators advocated violence.

0.67

0.33

# Unsupervised Approach

$$\text{argmax}_i \, P_{LM}(s_1, s_2)$$

$s_1$: **The city councilmen refused the demonstrators a permit because the city councilmen advocated violence.**

$s_2$: **The city councilmen refused the demonstrators a permit because the demonstrators advocated violence.**

# Unsupervised Approach

$$\text{argmax}_i \, P_{LM}(s_1, s_2)$$

$s_1$: **The city councilmen refused the demonstrators a permit because the city councilmen advocated violence.**

$s_2$: **The city councilmen refused the demonstrators a permit because the demonstrators advocated violence.**

$$\text{argmax}_i \sum_j P_{LM_j}(s_1, s_2)$$

Katrina had the financial means to afford a new car while Monica did not, since ____ had a high paying job.

**Sentence:**

Katrina had the financial means to afford a new car while Monica did not, since [MASK] had a high paying job.

**Predictions:**

11.8%  ↵

8.8%  She

6.3%  I

6.2%  So

5.2%  Monica

←  Undo

job —is used for→ make money

high-paying job ⇢type of⇢ job

make money ←requires— spend money

spend money ←requires— buy

buy —something to→ car

car —is capable of→ cost a lot of money

http://conceptnet5.media.mit.edu/

job —is used for→ make money

make a lot of money —entails→ make money

high-paying job —type of→ job

high-paying job —is used for→ make a lot of money

buy —requires→ spend money

spend money —requires→ make money

spend a lot of money —entails→ spend money

spend a lot of money —requires→ make a lot of money

buy something that costs a lot of money —entails→ buy

buy something that costs a lot of money —requires→ spend a lot of money

buy —something to→ car

buy something that costs a lot of money —something to→ car

car —is capable of→ cost a lot of money

# Neurosymbolic Approach



| 0.57 | 0.43 |

Model

vector representation

job — is used for → make money ← entails — make a lot of money

high-paying job — type of → job

high-paying job — is used for → make a lot of money

make money — requires — spend money

make a lot of money — requires — spend a lot of money

spend money ← entails — spend a lot of money

spend money — requires — buy

spend a lot of money — requires — buy something that costs a lot of money

buy ← entails — buy something that costs a lot of money

buy something that costs a lot of money — something to → car

car — is capable of → cost a lot of money

Katrina had the financial means to afford a new car while Monica did not, since ____ had a high paying job.

**Incorporating External Knowledge into Neural Models**

# Recipe

# Incorporating External Knowledge into Neural Models

# Recipe

## Knowledge Source

Knowledge bases, extracted from text, hand-crafted rules

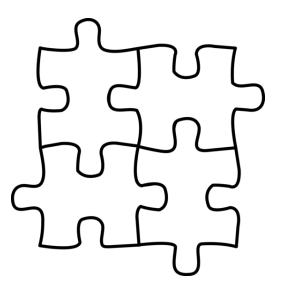# Incorporating External Knowledge into Neural Models

# Recipe

## Knowledge Source

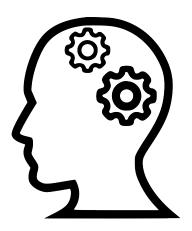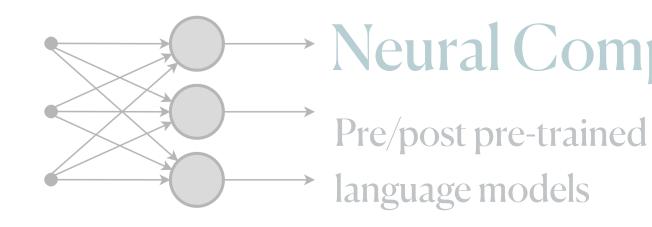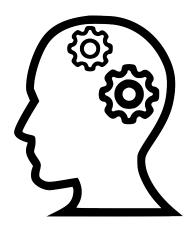Knowledge bases, extracted from text, hand-crafted rules

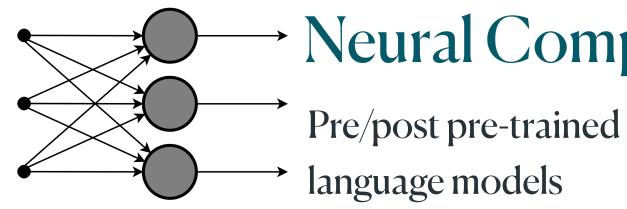## Neural Component

Pre/post pre-trained language models

# Incorporating External Knowledge into Neural Models
# Recipe

## Knowledge Source

Knowledge bases, extracted from text, hand-crafted rules

## Neural Component

Pre/post pre-trained language models

## Combination Method

Attention, pruning, word embeddings, multi-task learning

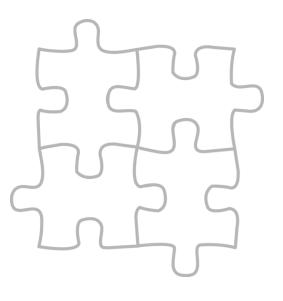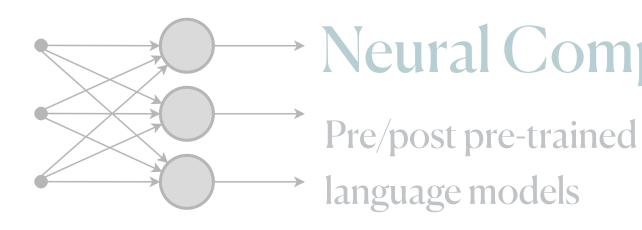# Incorporating External Knowledge into Neural Models

# Recipe

## Knowledge Source

Knowledge bases, extracted from text, hand-crafted rules

## Neural Component

Pre/post pre-trained language models

## Combination Method

Attention, pruning, word embeddings, multi-task learning

# Incorporating External Knowledge into Neural Models

# Recipe

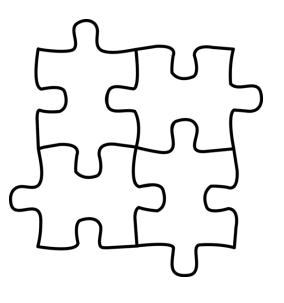## Knowledge Source

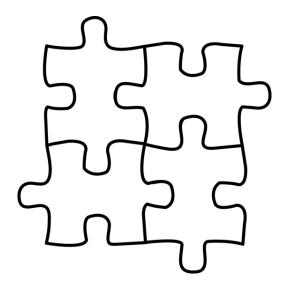Knowledge bases, extracted from text, hand-crafted rules

## Neural Component

Pre/post pre-trained language models

## Combination Method

Attention, pruning, word embeddings, multi-task learning

# Incorporating External Knowledge into Neural Models

# Recipe

### Knowledge Source
Knowledge bases, extracted from text, hand-crafted rules

### Neural Component
Pre/post pre-trained language models

### Combination Method
Attention, pruning, word embeddings, multi-task learning

# Combination Method

* Incorporate into scoring function
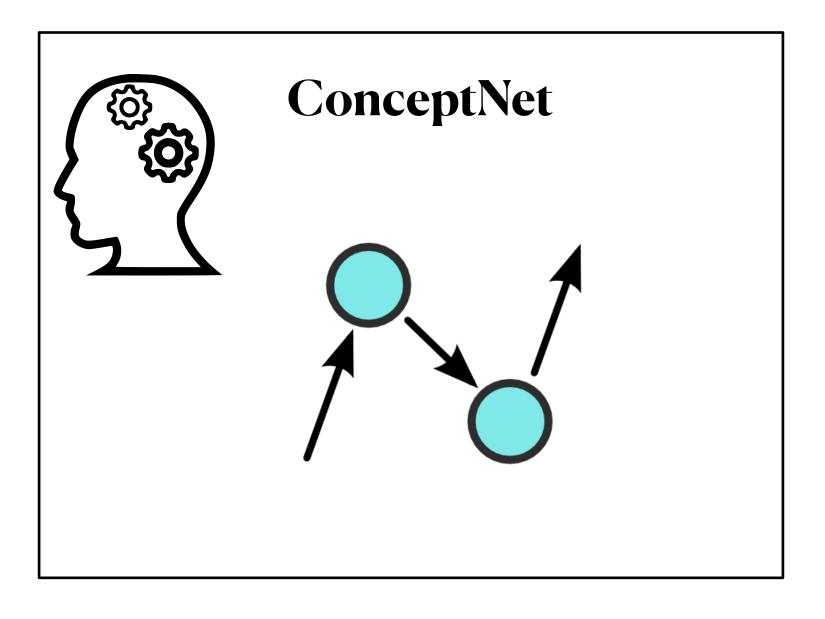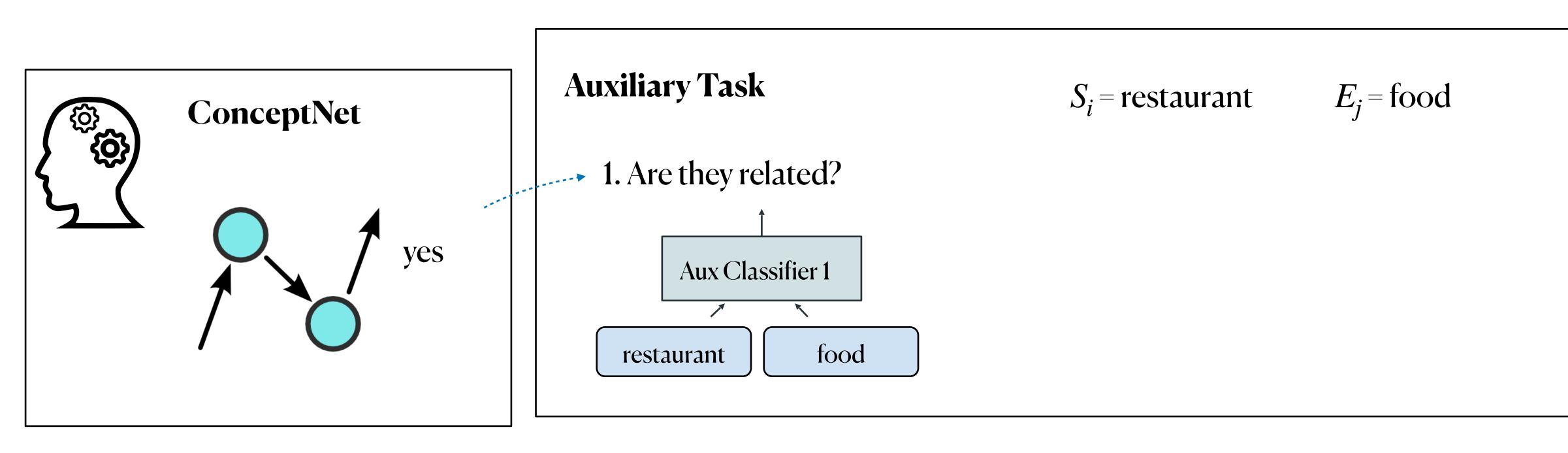* Multi-task learning
* Symbolic → vector representation (+attention)

# Combination Method

* Incorporate into scoring function
* **Multi-task learning**
* Symbolic → vector representation (+attention)

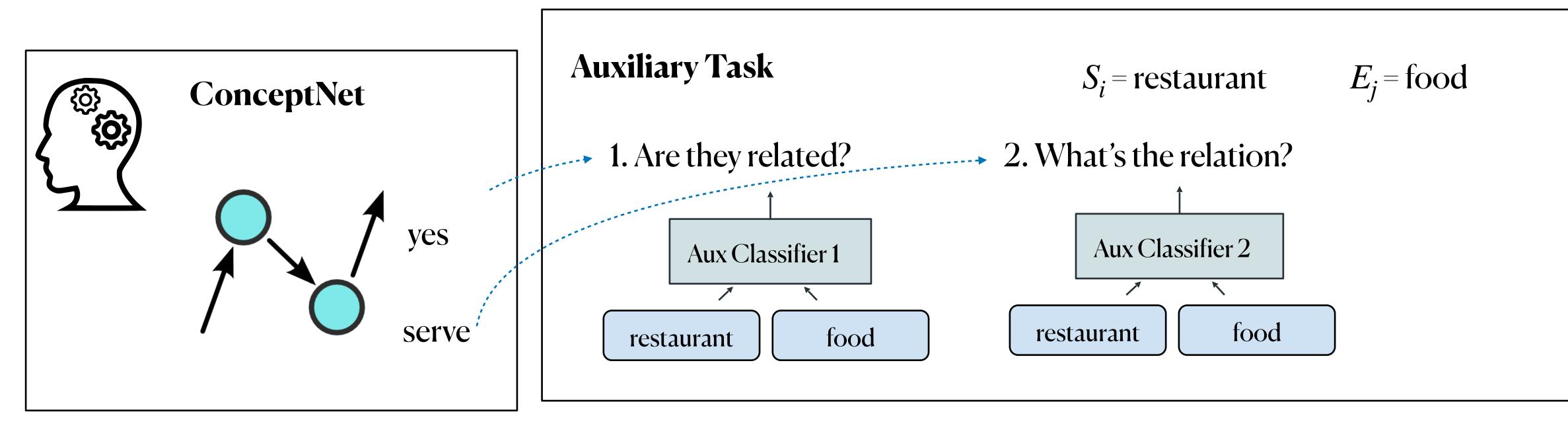# Incorporating External Knowledge into Neural Models
# Multitask Learning

Incorporating Relation Knowledge into Commonsense Reading Comprehension with Multi-task Learning. *Jiangnan Xia, Chen Wu, and Ming Yan.* CIKM 2019.

# Incorporating External Knowledge into Neural Models
# Multitask Learning



Incorporating Relation Knowledge into Commonsense Reading Comprehension with Multi-task Learning. *Jiangnan Xia, Chen Wu, and Ming Yan.* CIKM 2019.

# Incorporating External Knowledge into Neural Models
# Multitask Learning

# Incorporating External Knowledge into Neural Models
# Multitask Learning



ConceptNet

# Incorporating External Knowledge into Neural Models

# Multitask Learning



**ConceptNet**

yes

**Auxiliary Task**

$S_i$ = restaurant        $E_j$ = food

1. Are they related?

Aux Classifier 1

restaurant        food

# Incorporating External Knowledge into Neural Models

# Multitask Learning

# Limitations of Neurosymbolic Methods

# Limitations of Neurosymbolic Methods

‣ Knowledge graphs have **limited coverage**

Commonsense knowledge is **immeasurably vast**, making it **impossible to manually enumerate**

# Limitations of Neurosymbolic Methods

‣ Knowledge graphs have **limited coverage**

‣ Inferences may be correct only in certain **contexts**



en **mouse**
An English term in ConceptNet 5.8

**Sources:** Open Mind Common Sense contributors, DBPedia 201!
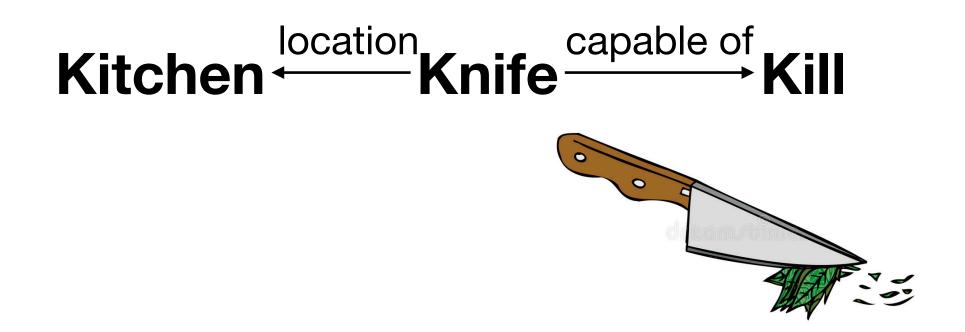WordNet
View this term in the API

**Location of mouse**

en a hole in a wall →
en the garage →
en a laboratory →
en the attic →
en a cupboard →
en a kitchen →
en a trap →
en a cellar →
en your desk →
en a hole →
en sewer →

# Limitations of Neurosymbolic Methods

‣ Knowledge graphs have **limited coverage**

‣ Inferences may be correct only in certain **contexts**

‣ Long KB paths have **limited precision**

**Kitchen** ←——location——— **Knife** ——capable of——→ **Kill**

# Limitations of Neurosymbolic Methods

‣ Knowledge graphs have **limited coverage**

‣ Inferences may be correct only in certain **contexts**

‣ Long KB paths have **limited precision**

‣ Tradeoff: embedding knowledge (**better generalization**)
  vs. hard constraints (**more accurate**)

# Limitations of Neurosymbolic Methods

‣ Knowledge graphs have **limited coverage**

‣ Inferences may be correct only in certain **contexts**

‣ Long KB paths have **limited precision**

‣ Tradeoff: embedding knowledge (**better generalization**) vs. hard constraints (**more accurate**)

# COMET

Given a **seed entity** and a **relation**,
learn to generate the **target entity**

**tail entity**

Language Model

↑ person    ↑ sails    ↑ across    ↑ oceans    ↑ <requires>

_____  _____
**head entity**              **relation**

COMET: Commonsense Transformers for Automatic Knowledge Graph Construction. Antoine Bosselut, Hannah Rashkin, Maarten Sap, Chaitanya Malaviya, Asli Celikyilmaz, and Yejin Choi. ACL 2020

# COMET

Given a **seed entity** and a **relation**,
learn to generate the **target entity**

tail entity

buy

## Language Model

person   sails   across   oceans   &lt;requires&gt;

head entity

relation

# COMET

Given a **seed entity** and a **relation**, learn to generate the **target entity**

tail entity

buy          a

## Language Model

person    sails    across    oceans    &lt;requires&gt;

head entity          relation

COMET: Commonsense Transformers for Automatic Knowledge Graph Construction. Antoine Bosselut, Hannah Rashkin, Maarten Sap, Chaitanya Malaviya, Asli Celikyilmaz, and Yejin Choi. ACL 2020
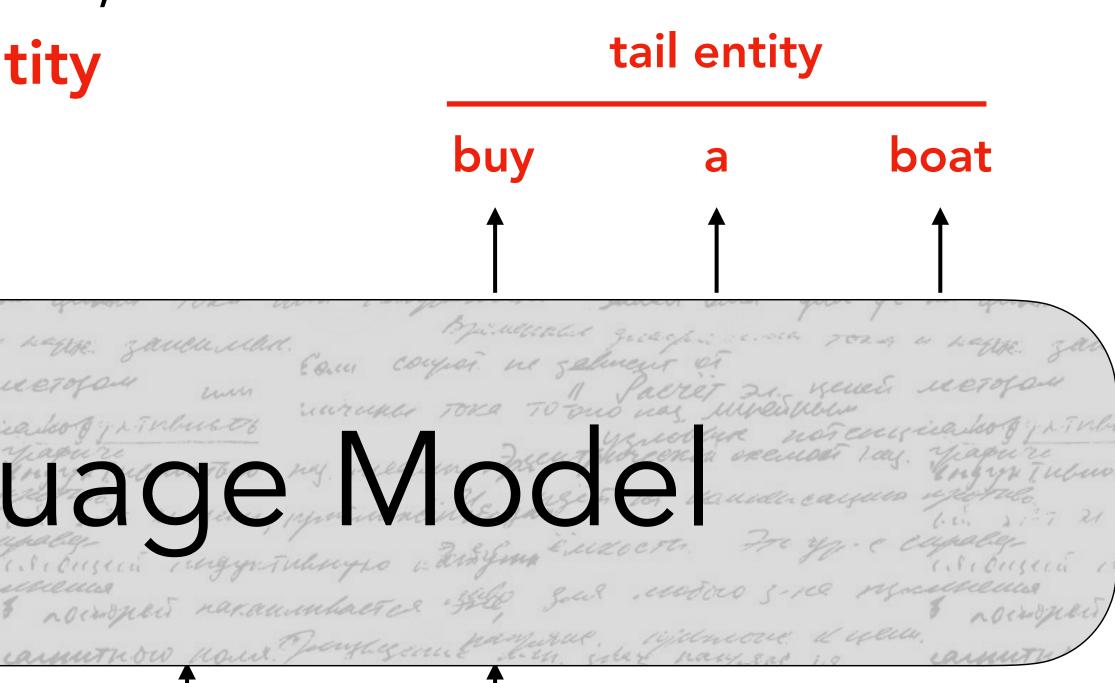
# COMET

Given a **seed entity** and a **relation**,
learn to generate the **target entity**

<u>tail entity</u>

buy        a        boat

Language Model

person    sails    across    oceans    &lt;requires&gt;

<u>head entity</u>         <u>relation</u>

COMET: Commonsense Transformers for Automatic Knowledge Graph Construction. Antoine Bosselut, Hannah Rashkin, Maarten Sap, Chaitanya Malaviya, Asli Celikyilmaz, and Yejin Choi. ACL 2020
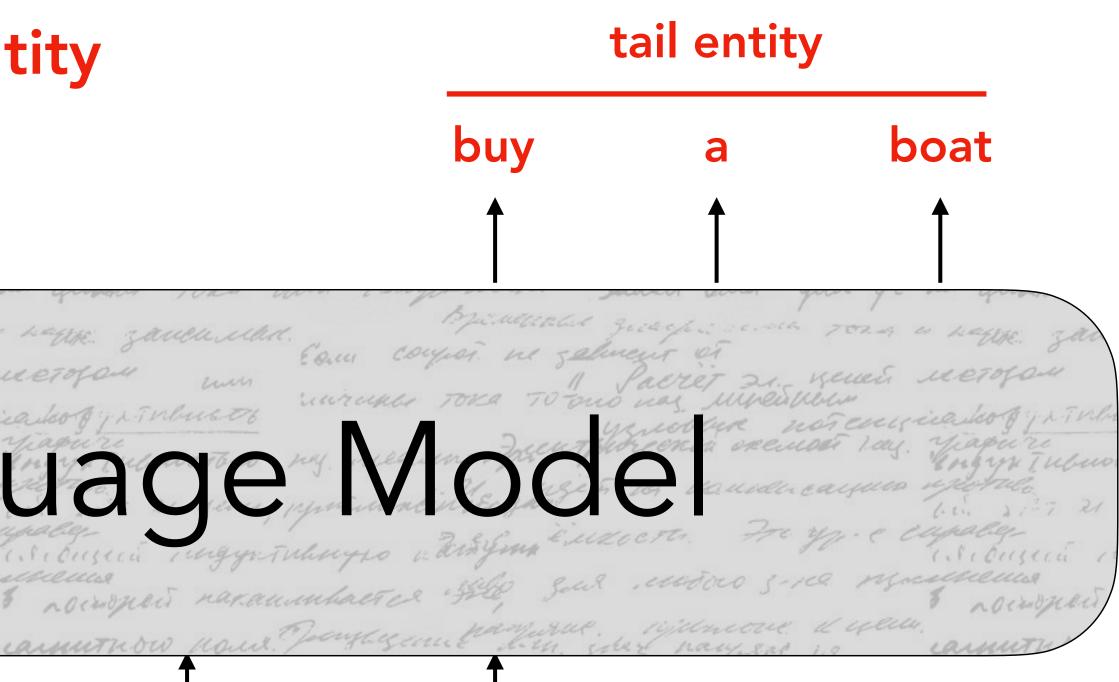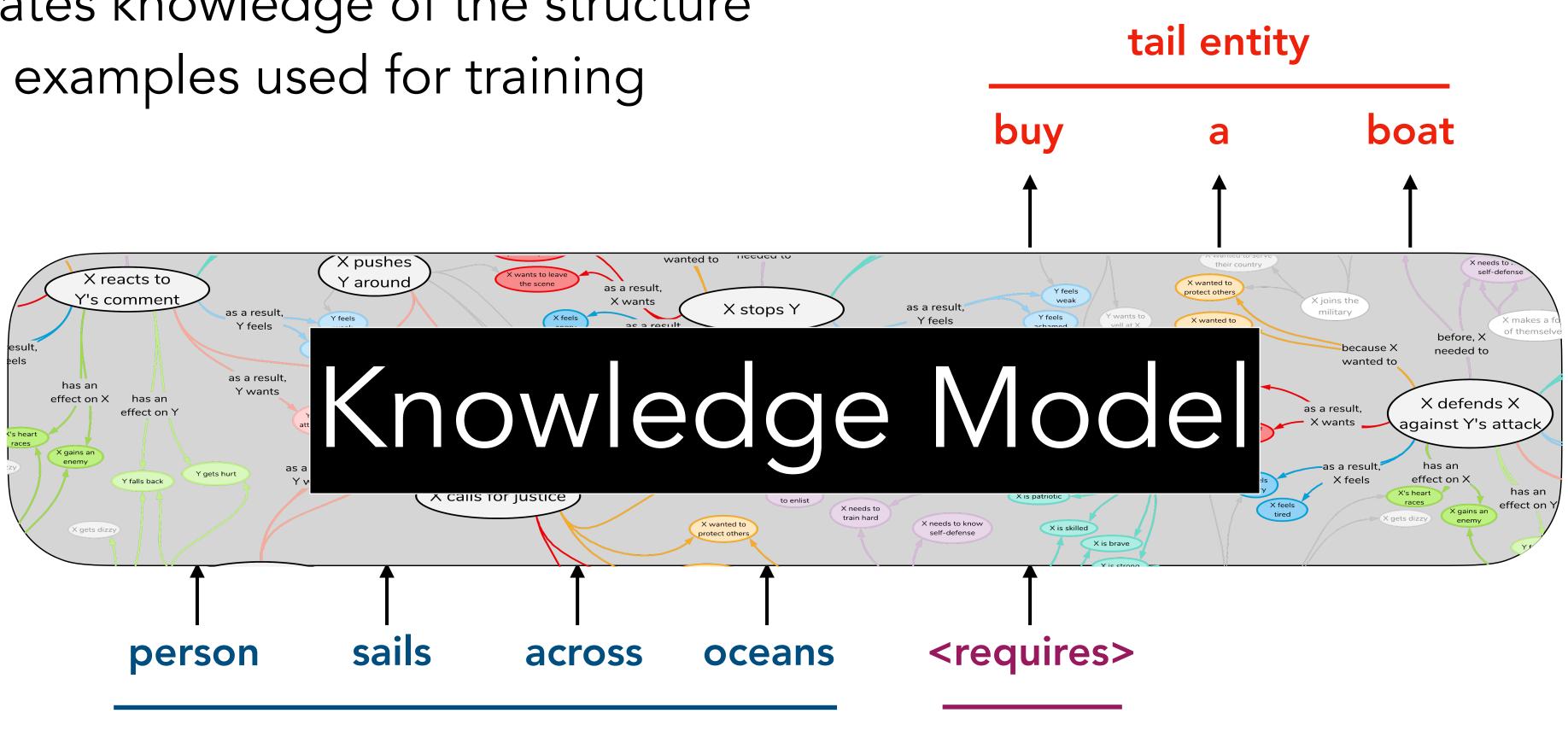
# COMET

Given a **seed entity** and a **relation**, learn to generate the **target entity**

$$\mathcal{L} = -\sum \log P(\textbf{target words} \mid \textbf{seed words}, \textbf{relation})$$

**tail entity**

**buy**     **a**     **boat**



Language Model

person     sails     across     oceans     \<requires\>

**head entity**     **relation**

COMET: Commonsense Transformers for Automatic Knowledge Graph Construction. Antoine Bosselut, Hannah Rashkin, Maarten Sap, Chaitanya Malaviya, Asli Celikyilmaz, and Yejin Choi. ACL 2020
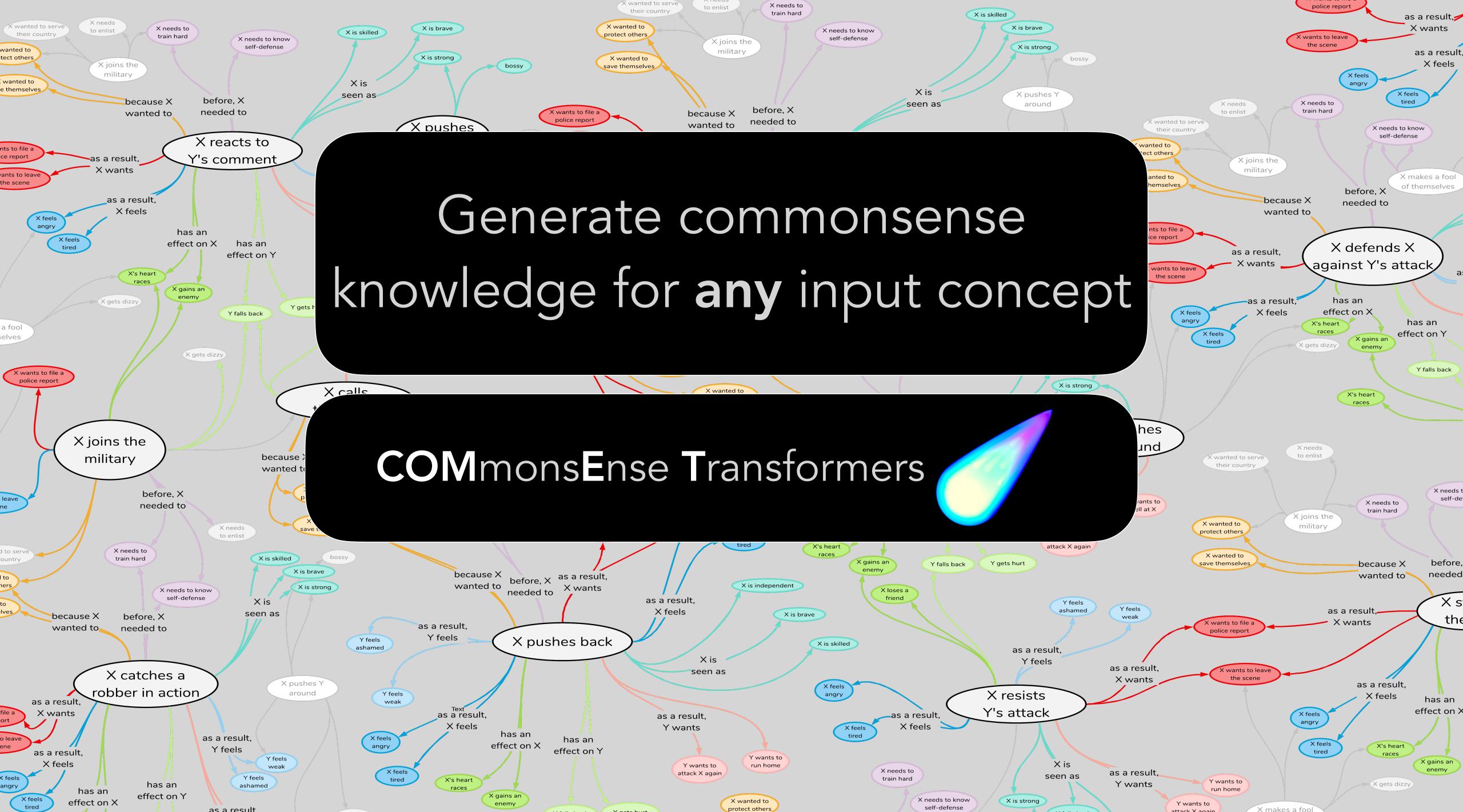
# COMET

Language Model → **Knowledge Model:**
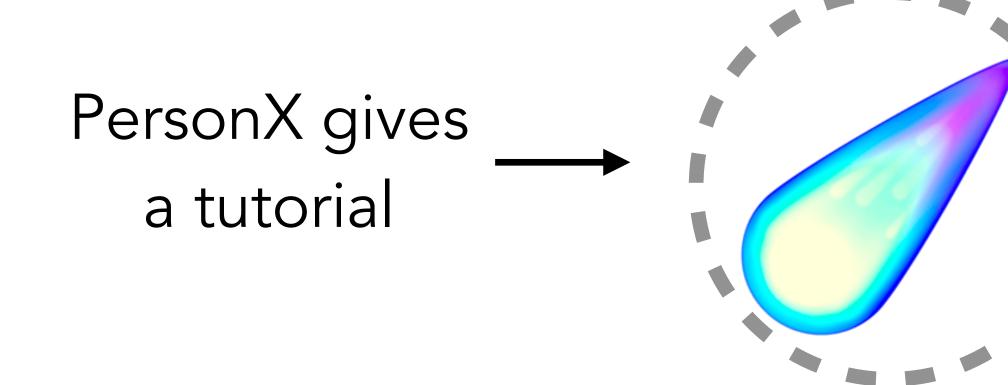generates knowledge of the structure
of the examples used for training

COMET: Commonsense Transformers for Automatic Knowledge Graph Construction. Antoine Bosselut, Hannah Rashkin, Maarten Sap, Chaitanya Malaviya, Asli Celikyilmaz, and Yejin Choi. ACL 2020
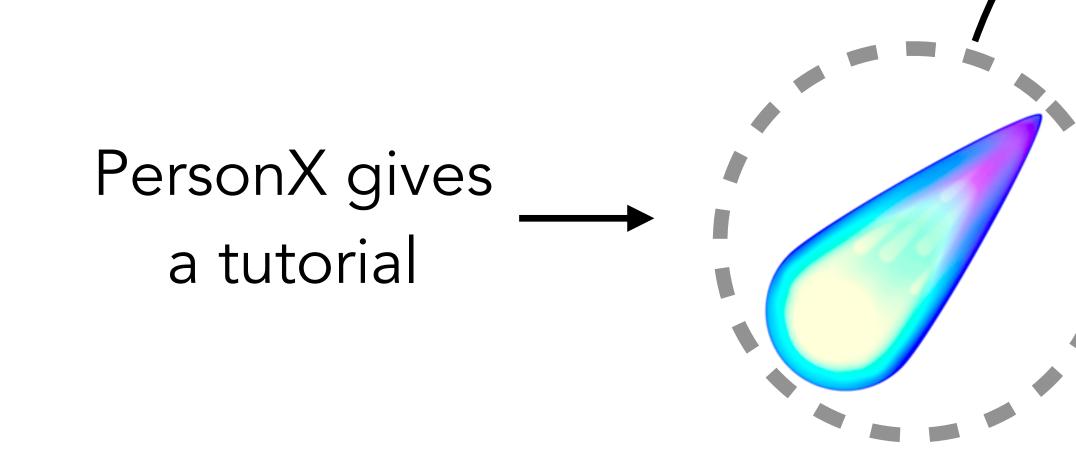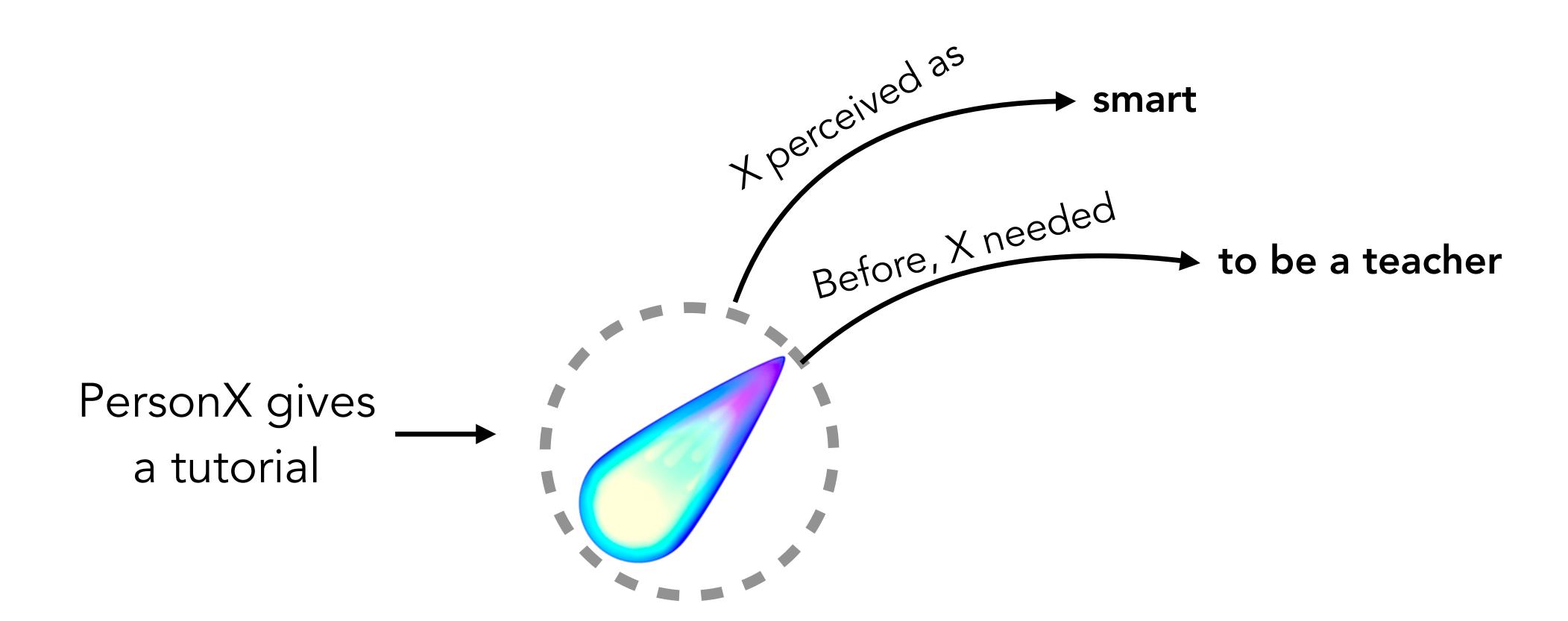
Generate commonsense knowledge for **any** input concept

Generate commonsense knowledge for **any** input concept

**COM**mons**E**nse **T**ransformers

# COMET - ATOMIC

PersonX gives
a tutorial →

# COMET - ATOMIC

# COMET - ATOMIC



X perceived as → **smart**

Before, X needed → **to be a teacher**

PersonX gives a tutorial

# COMET - ATOMIC



PersonX gives a tutorial

X perceived as → **smart**

Before, X needed → **to be a teacher**

Others will want → **to thank PersonX**

# COMET - ATOMIC
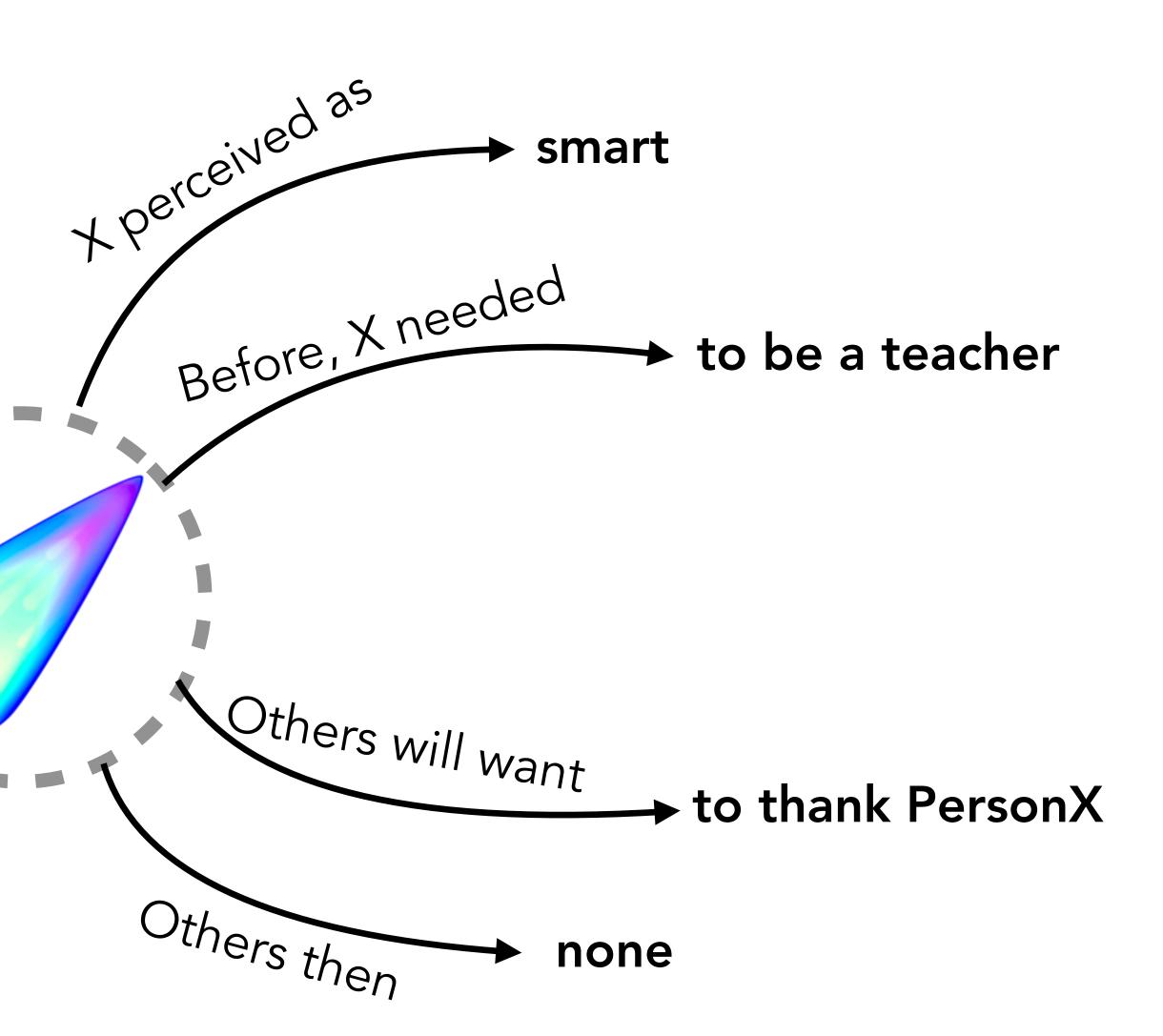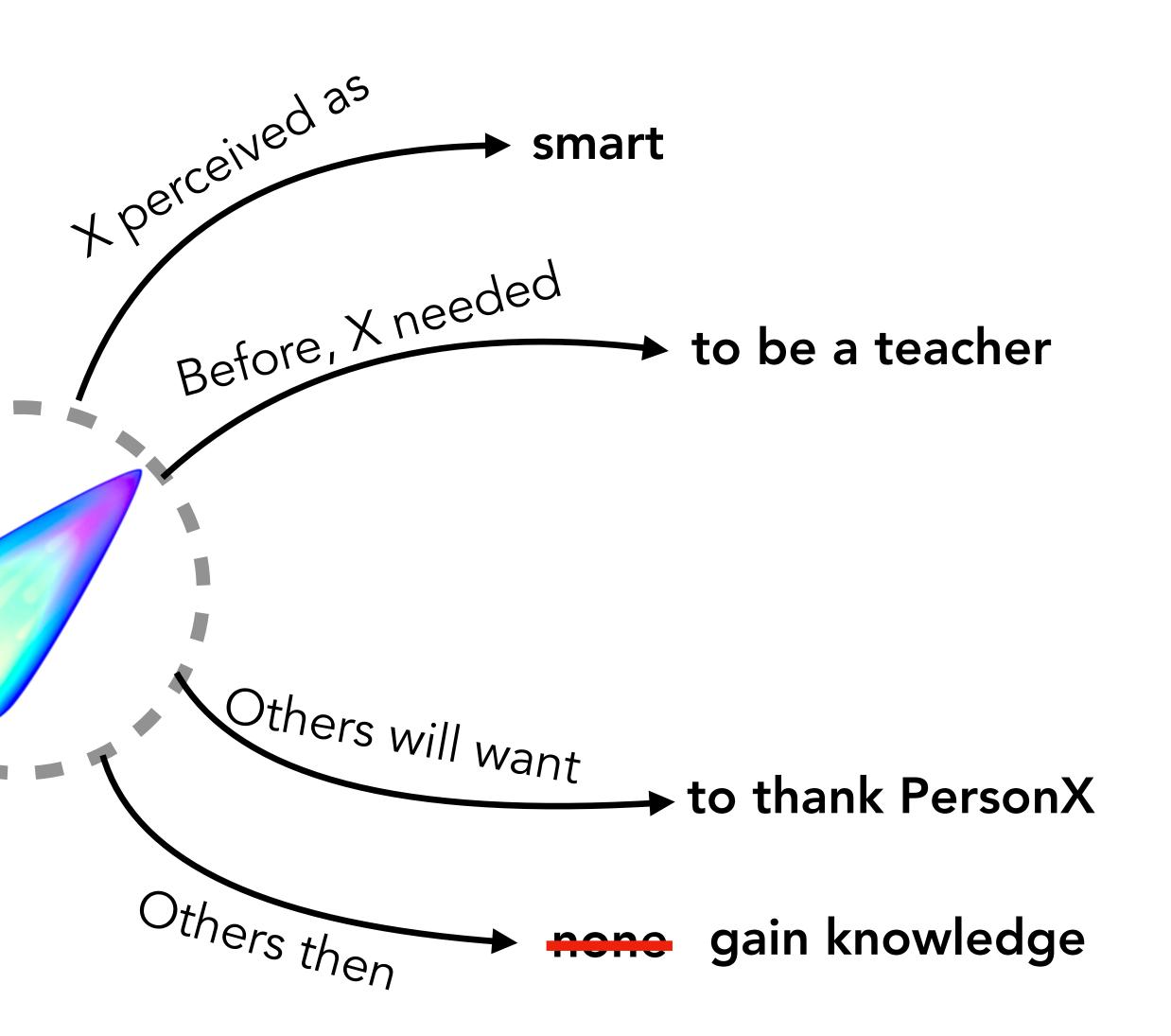


PersonX gives a tutorial

X perceived as → **smart**

Before, X needed → **to be a teacher**

Others will want → **to thank PersonX**

Others then → **none**

# COMET - ATOMIC
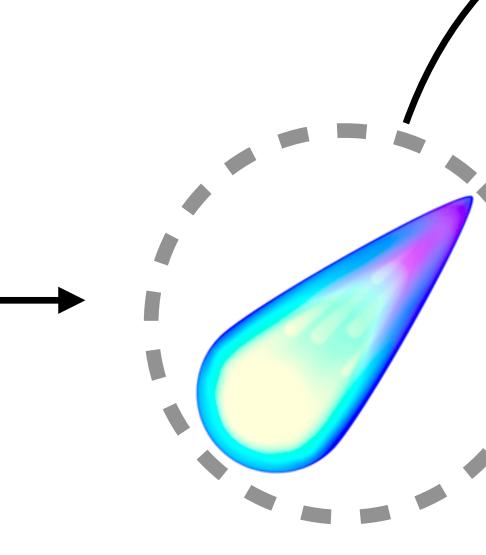
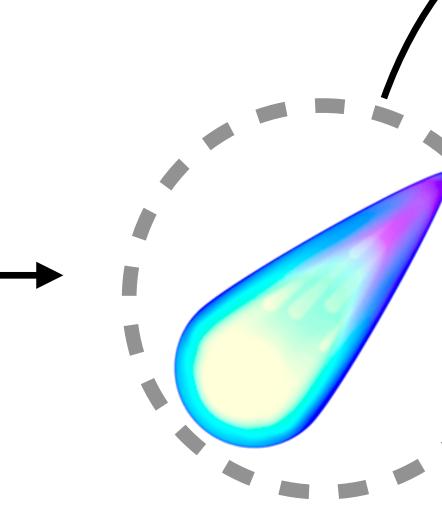

PersonX gives a tutorial

X perceived as → **smart**

Before, X needed → **to be a teacher**

Others will want → **to thank PersonX**

Others then → ~~none~~ **gain knowledge**

# COMET - ConceptNet

listen to
tutorial →

# COMET - ConceptNet



listen to tutorial →

location → **classroom**

# COMET - ConceptNet

listen to tutorial →

location → **classroom**

motivated by → **you be smart**

# COMET - ConceptNet



listen to tutorial →

location → **classroom**

motivated by → **you be smart**

starts with → **sit down**

# COMET - ConceptNet



listen to tutorial →

location → **classroom**

motivated by → **you be smart**

starts with → **sit down**

has prerequisite → **listen carefully**

# COMET - ConceptNet



listen to tutorial →

location → **classroom**

motivated by → **you be smart**

starts with → **sit down**

has prerequisite → **listen carefully**

causes → **good grade**

# Recap

# Recap

**Benchmarks:**

- ☑ Measure progress

- ☑ Cover different types of knowledge & reasoning

- ☐ Tradeoff:
  easy to evaluate vs.
  hard to game

# Recap

**Benchmarks:**

☑ Measure progress

☑ Cover different types of knowledge & reasoning

☐ Tradeoff:
easy to evaluate vs.
hard to game

**Symbolic Knowledge:**

☑ Use for completing missing / unstated knowledge

☐ Insufficient coverage
☐ How to collect?
☐ How to incorporate into models?

# Recap

**Benchmarks:**
- ☑ Measure progress
- ☑ Cover different types of knowledge & reasoning
- ☐ Tradeoff:
  easy to evaluate vs.
  hard to game

**Symbolic Knowledge:**
- ☑ Use for completing missing / unstated knowledge
- ☐ Insufficient coverage
- ☐ How to collect?
- ☐ How to incorporate into models?

**Neural Representations:**
- ☑ Generalization
- ☑ Easy to train/use
- ☐ Inaccurate

# Recap

**Benchmarks:**
- ☑ Measure progress
- ☑ Cover different types of knowledge & reasoning
- ☐ Tradeoff: easy to evaluate vs. hard to game

**Symbolic Knowledge:**
- ☑ Use for completing missing / unstated knowledge
- ☐ Insufficient coverage
- ☐ How to collect?
- ☐ How to incorporate into models?

**Neural Representations:**
- ☑ Generalization
- ☑ Easy to train/use
- ☐ Inaccurate

**Reasoning Engine:**
- ☐ We're only scratching the surface... and we didn't really talk about "reasoning"!

# Recap

**Benchmarks:**
- ☑ Measure progress
- ☑ Cover different types of knowledge & reasoning
- ☐ Tradeoff: easy to evaluate vs. hard to game

**Symbolic Knowledge:**
- ☑ Use for completing missing / unstated knowledge
- ☐ Insufficient coverage
- ☐ How to collect?
- ☐ How to incorporate into models?

**Neural Representations:**
- ☑ Generalization
- ☑ Easy to train/use
- ☐ Inaccurate

**Reasoning Engine:**
- ☐ We're only scratching the surface... and we didn't really talk about "reasoning"!

*Thank You!*

@VeredShwartz          vshwartz@cs.ubc.ca