

# Symbolic Knowledge Distillation

From General Language Models to **Commonsense** Models

— NAACL 2022 —



Peter  
West

Chandra  
Bhagavatula



Jack  
Hessel



Jena  
Hwang



Liwei  
Jiang



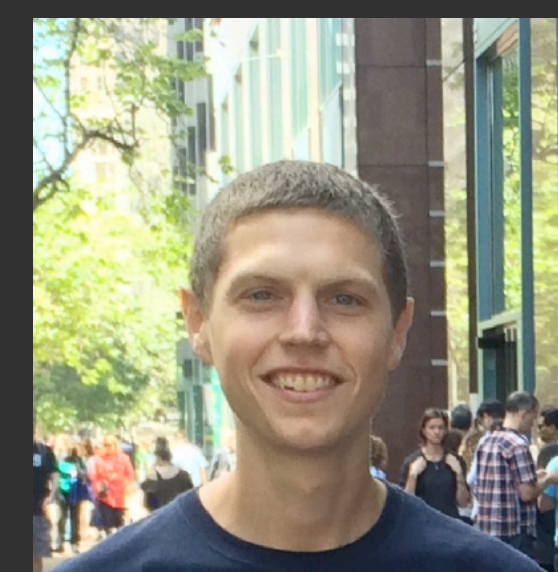
Ronan  
Le Bras



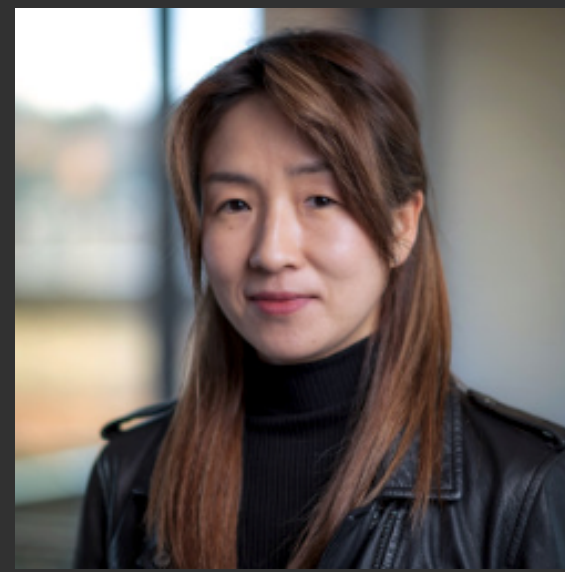
Ximing  
Lu



Sean  
Welleck



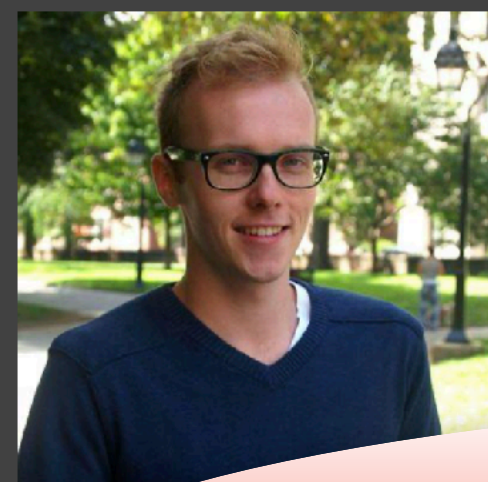
Yejin  
Choi



# Language models != knowledge models

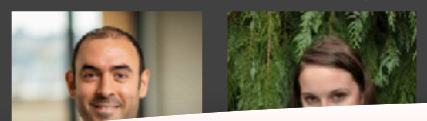
ATOMIC: An  
Com  
for If-Th

Maarten Sap



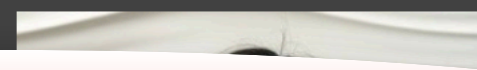
Ronan  
LeBras

Emily  
Allaway



Jena  
Hwang

Chandra  
Bhagavatula



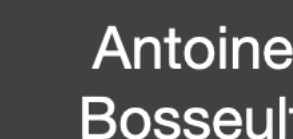
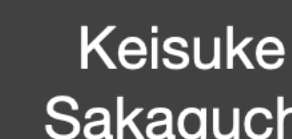
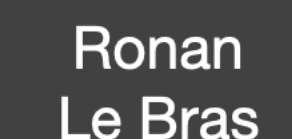
Ronan  
Le Bras

Jeff  
Da

Keisuke  
Sakaguchi

Antoine  
Bosseult

Me



Fully crowdsourced by humans

**Symbolic** commonsense  
knowledge graph

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

On Symbolic and Neural Commonsense Knowledge Graphs

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

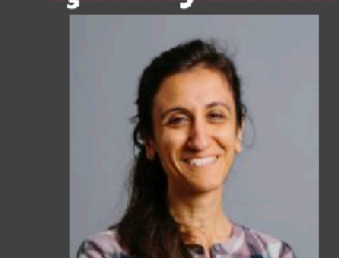
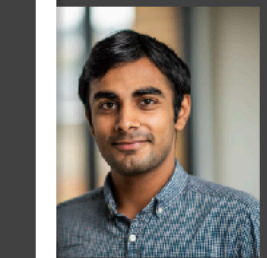
AAAI 2021

Transformers for  
Graph Construction

Shantanya  
Malaviya

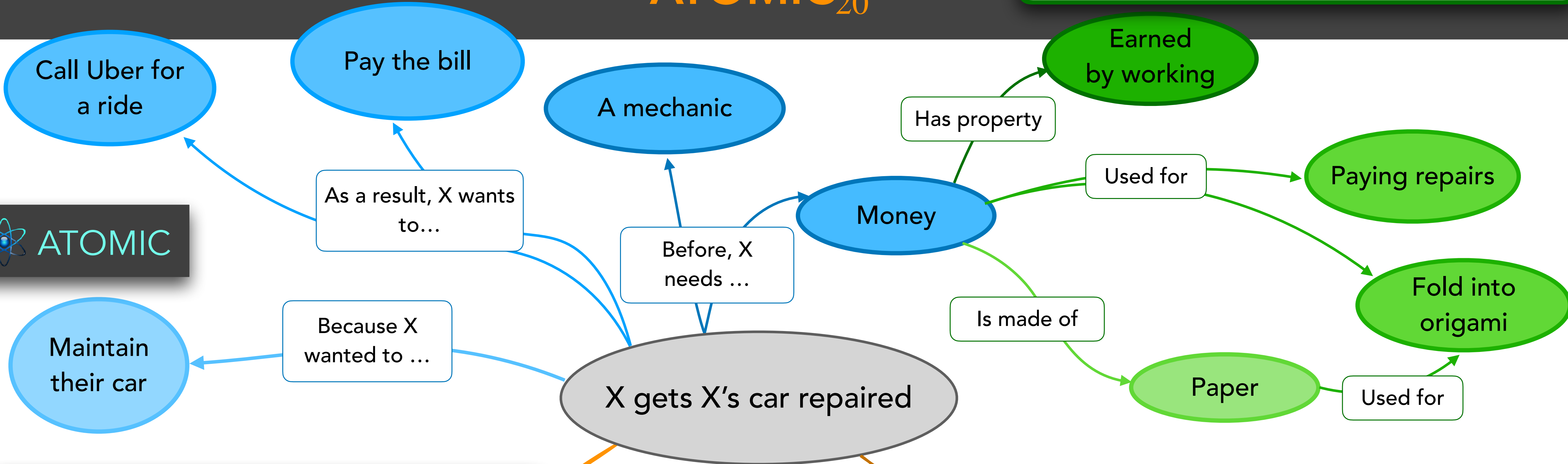
Asli  
Çelikyilmaz

Me

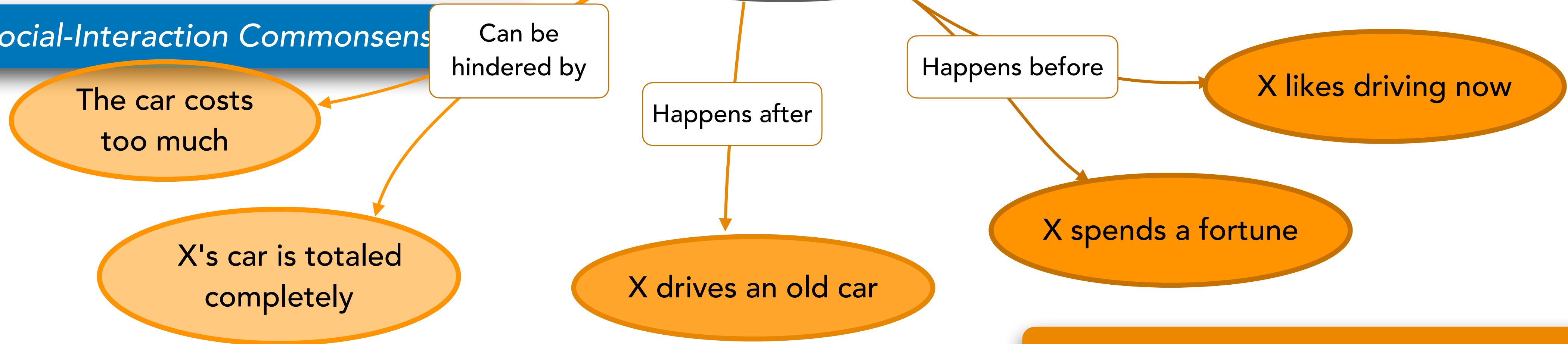


**Neural** commonsense model

## Physical-Entity Commonsense

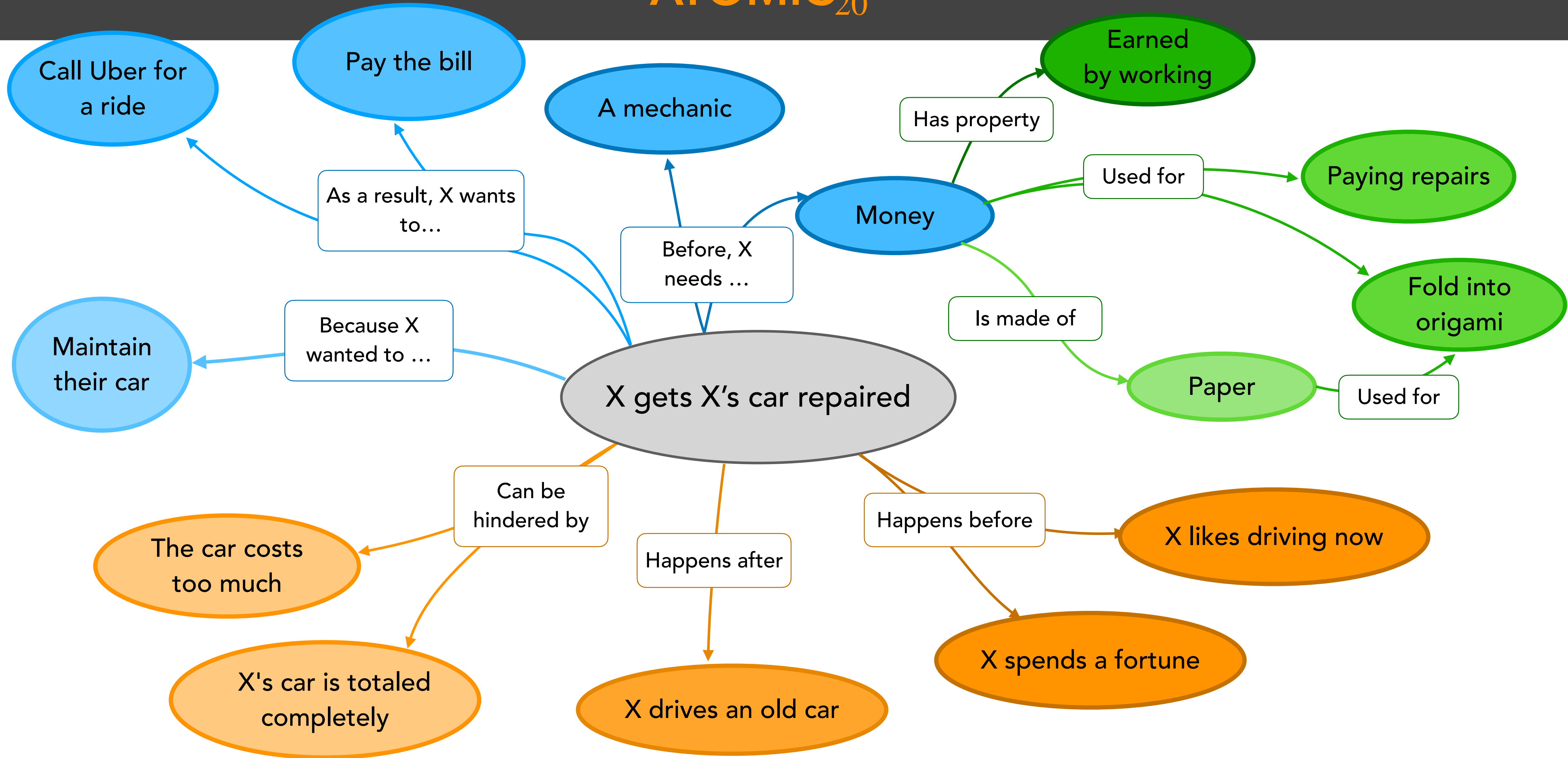


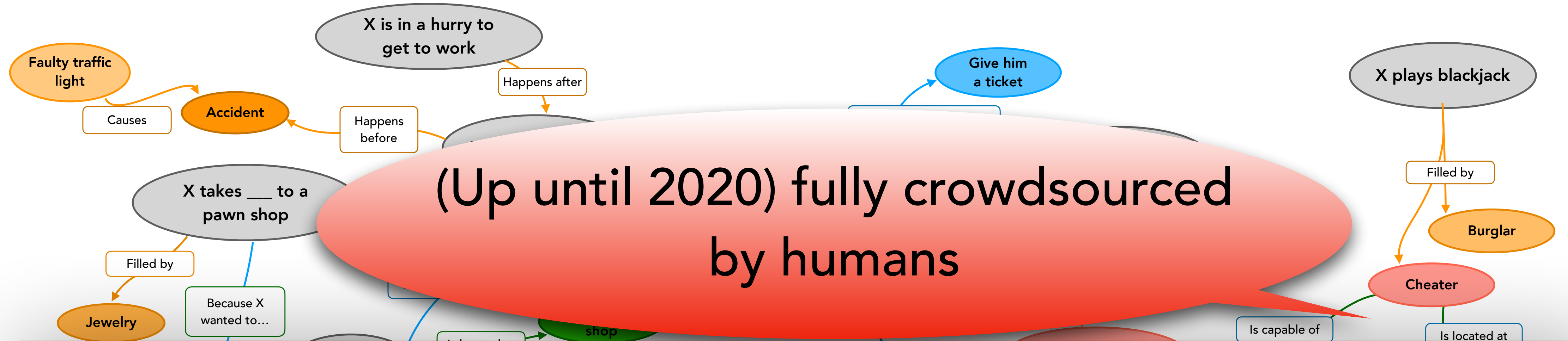
## Social-Interaction Commonsense



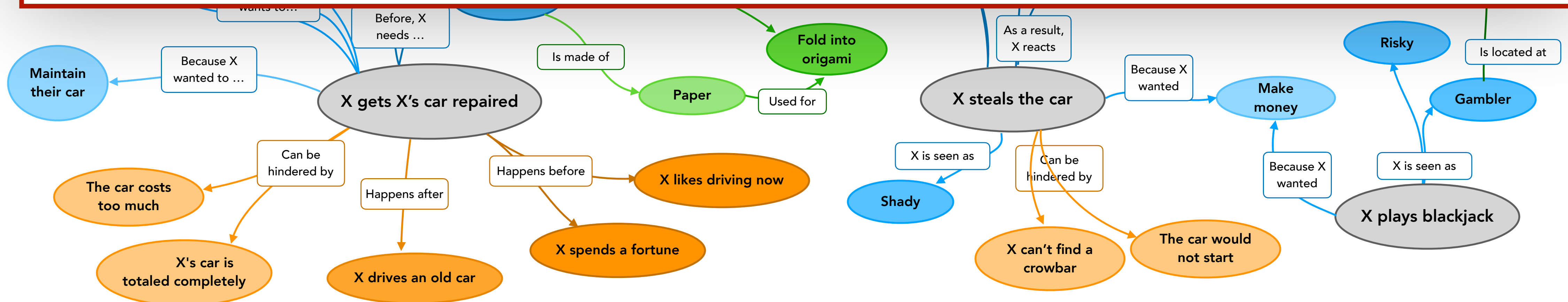
## Event-Centered Commonsense

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



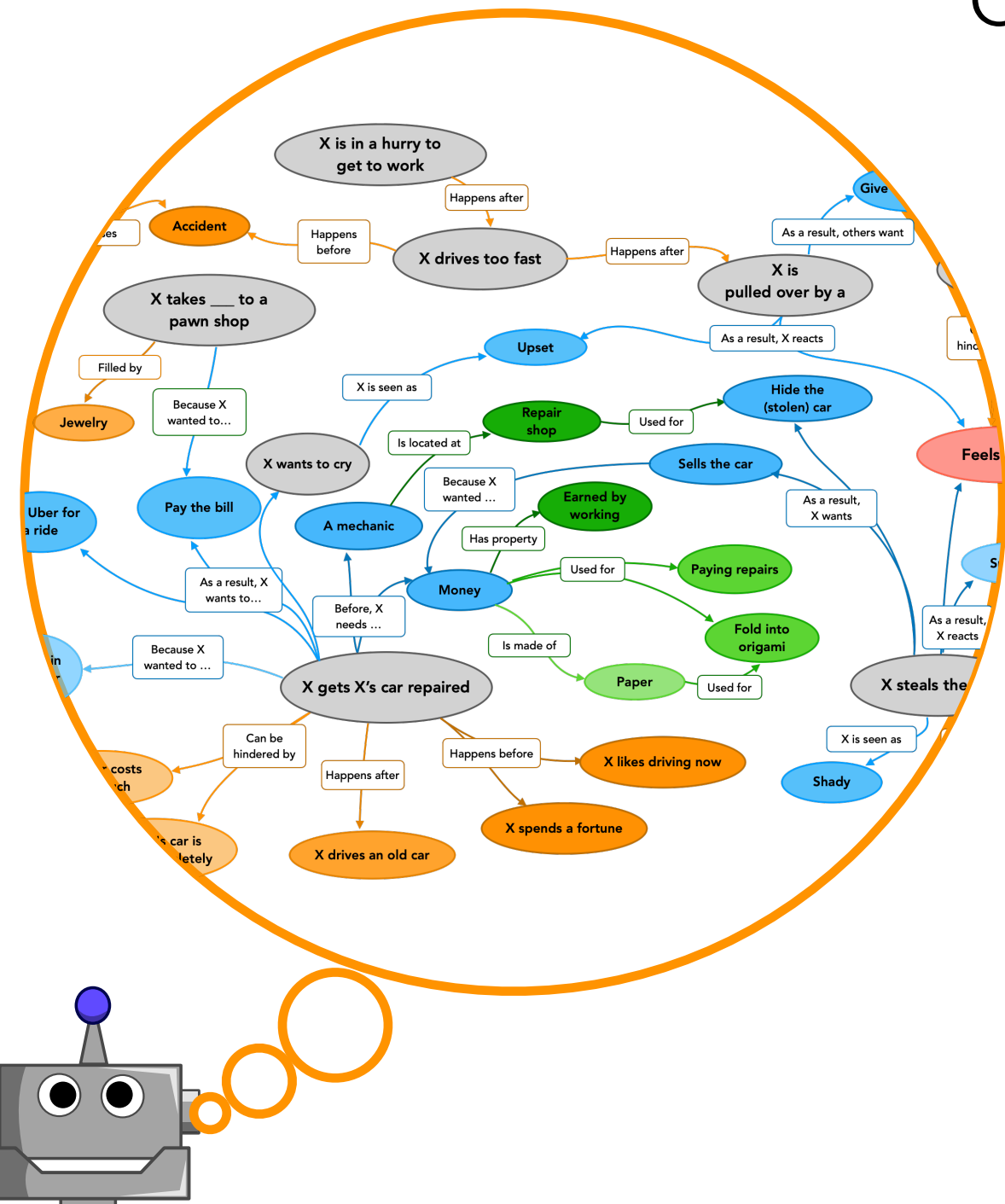
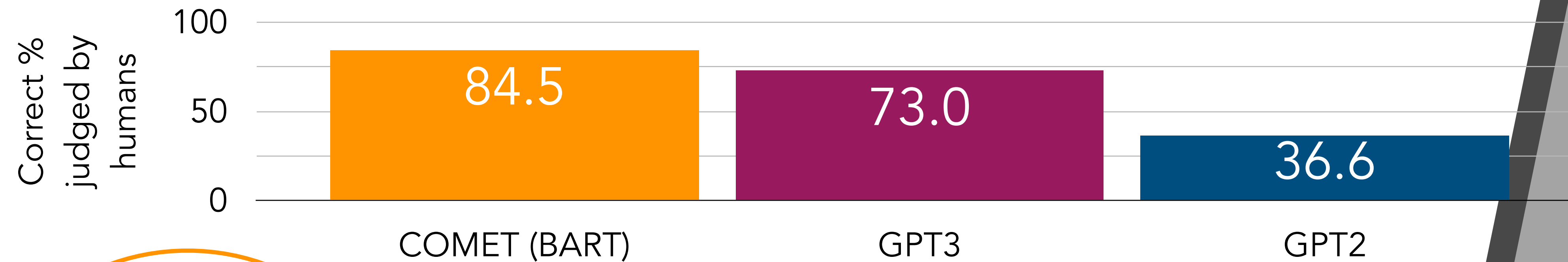


**1.33M commonsense if-then inferences**  
**23 relations (or inference types)**



## Knowledge Models

## Off-the-shelf Language Models



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

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

## Persona-aware Conversations

**Like Hiking? Person-grounded Dialog**  
(Majumder et al, 2020)  
EMNLP '20

**Health Counseling Dialogue**  
(Kearns et al, 2020)  
CHI EA '20

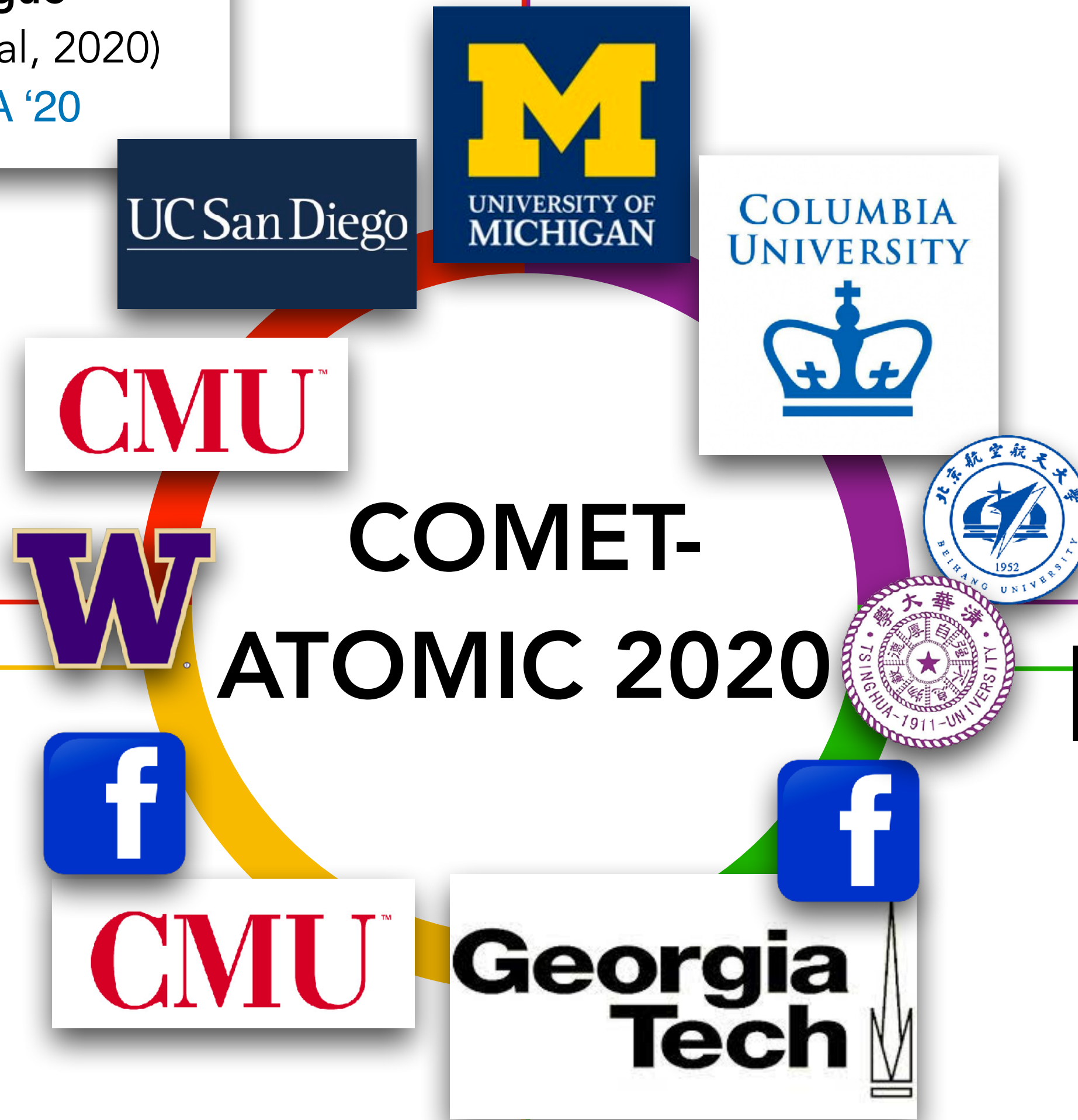
**COSMIC: Emotion Identification in Conversations**  
(Ghosal et al, 2020)  
EMNLP '20

## Figurative Language Understanding

**Metaphor Generation with Conceptual Mapping**  
(Stowe et al, 2021)  
ACL '21

**MERMAID: Metaphor Generation**  
(Chakrabarty et al, 2021)  
NAACL '21

# COMET-ATOMIC 2020



## Interactive Learning Enhancement

**Conversation Multi-hop Reasoning through Neural Commonsense**  
(Forough et al, 2021)  
EMNLP '21

## Storytelling and Fantasy Gaming

**How to Motivate Your Dragon**  
(Ammanabrolu et al, 2021)  
AAAI '21

**Commonsense Story Generation**  
(Guan et al, 2020)  
TACL '20

# Symbolic Knowledge Distillation

From Neural Language Models to **Causal Commonsense** Models

*New:*

*ATOMIC-10x*

*COMET-distill*



Peter  
West

Chandra  
Bhagavatula



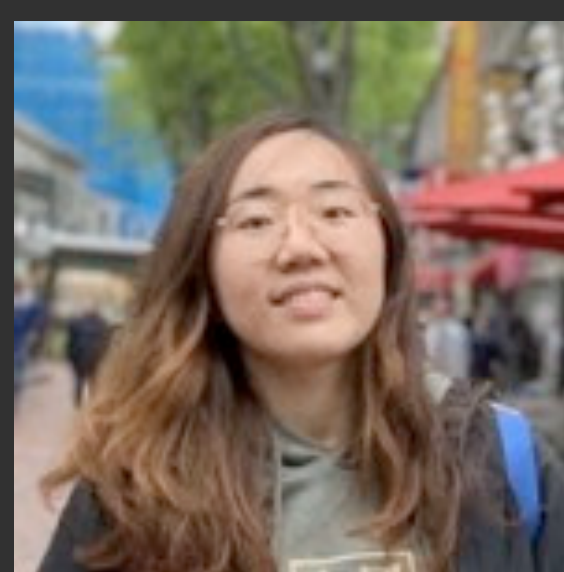
Jack  
Hessel



Jena  
Hwang



Liwei  
Jiang



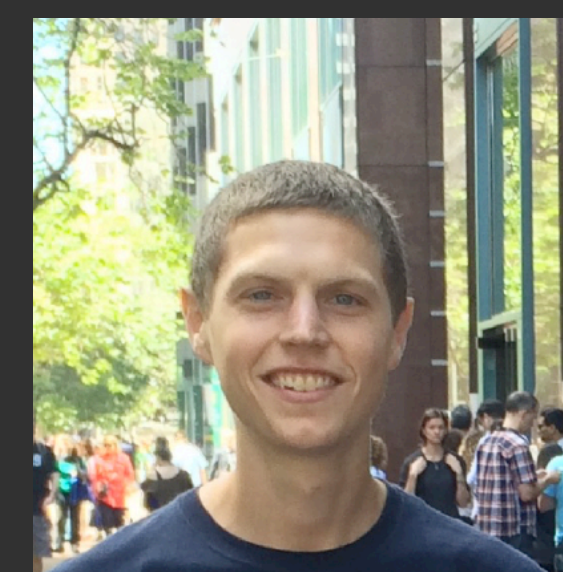
Ronan  
Le Bras



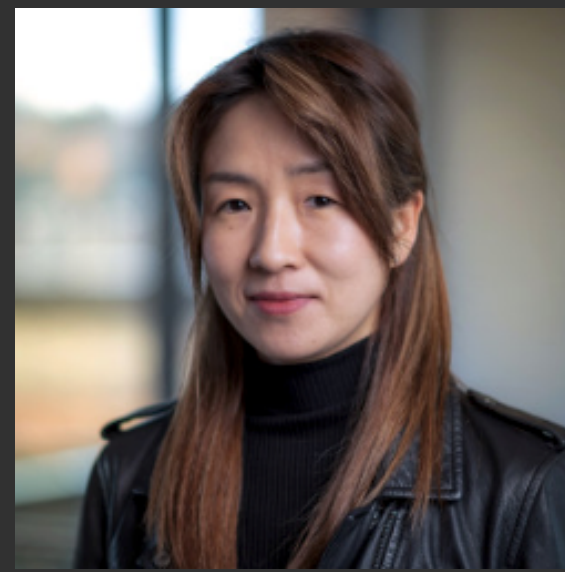
Ximing  
Lu



Sean  
Welleck

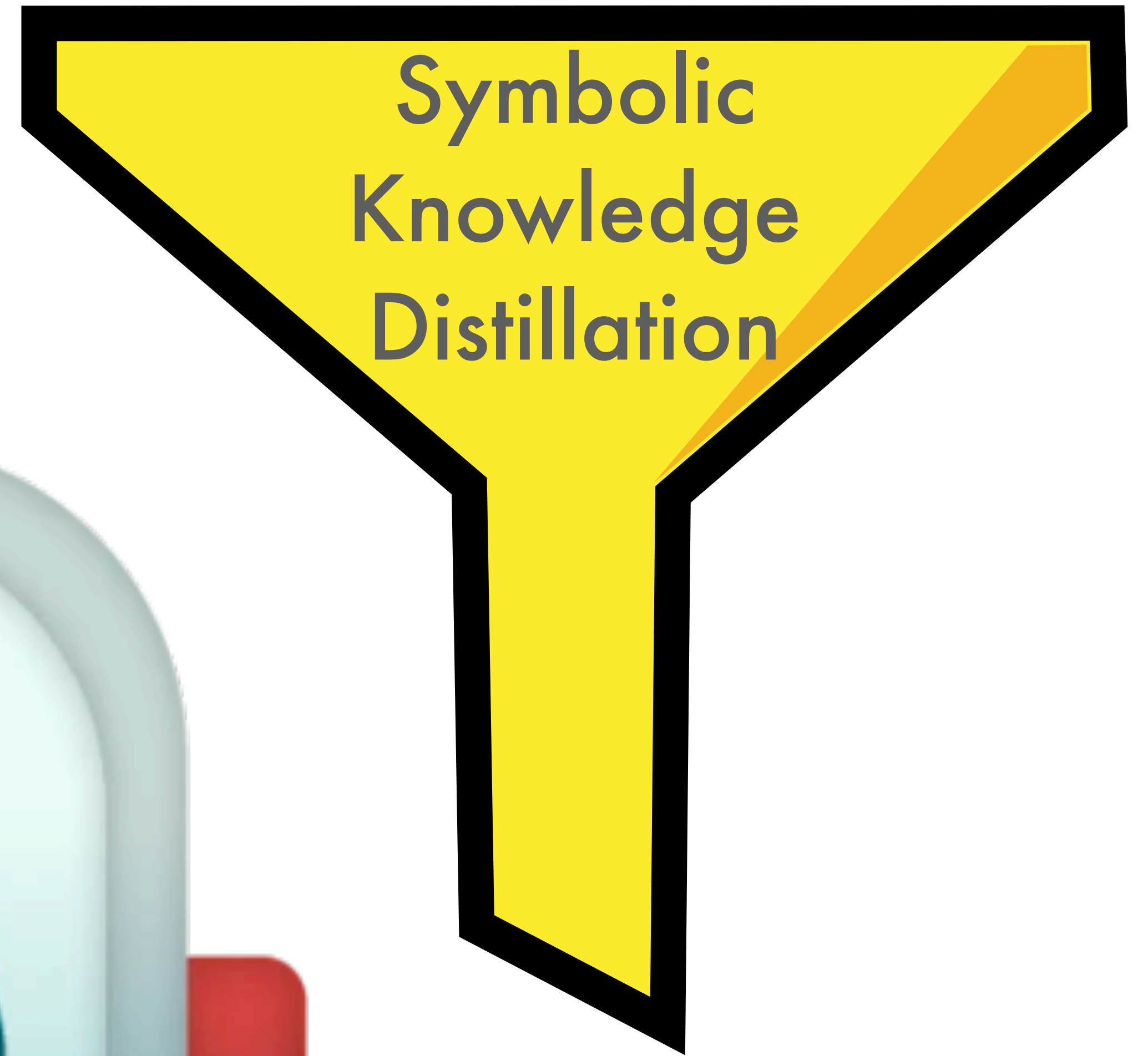


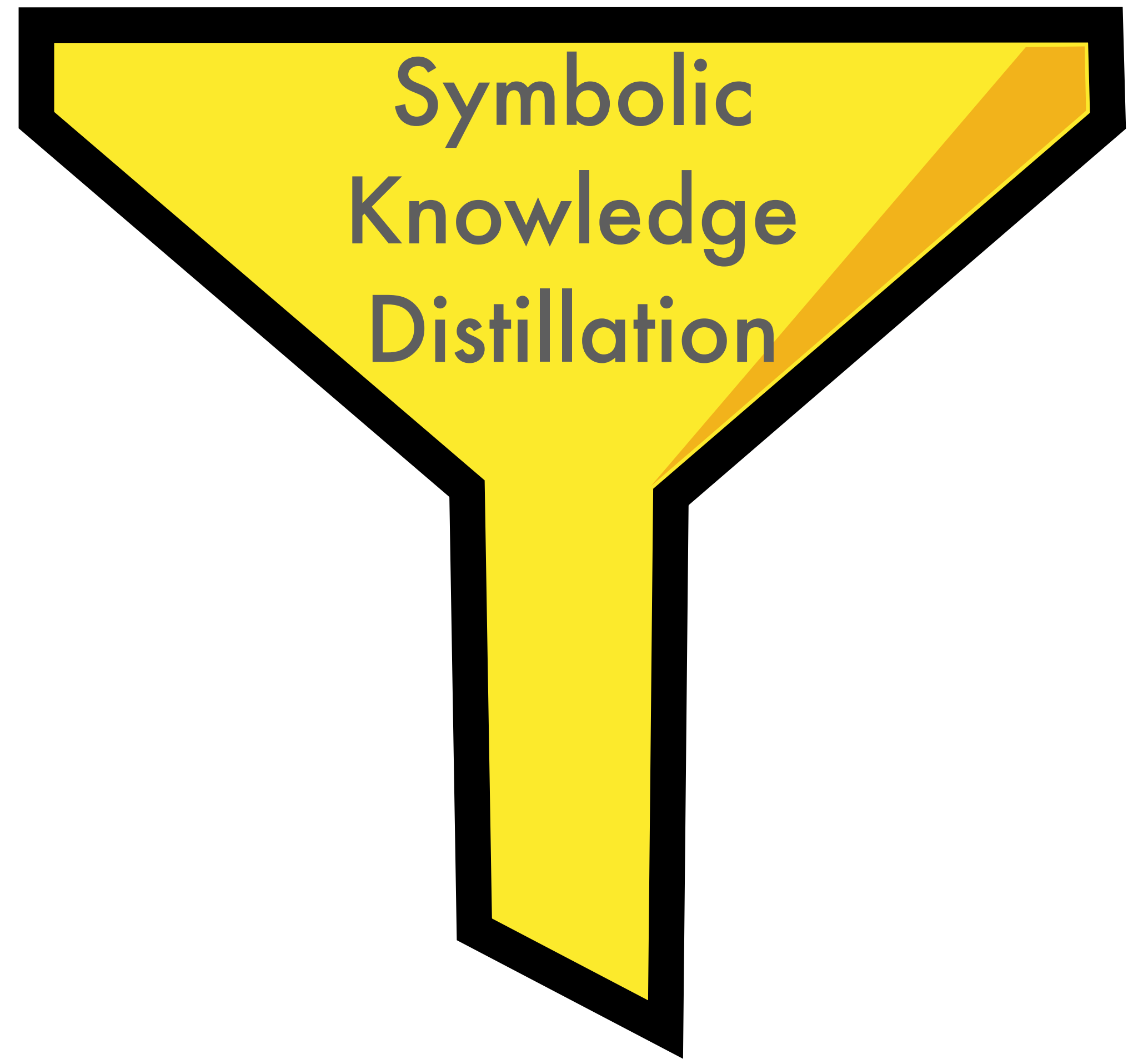
Yejin  
Choi

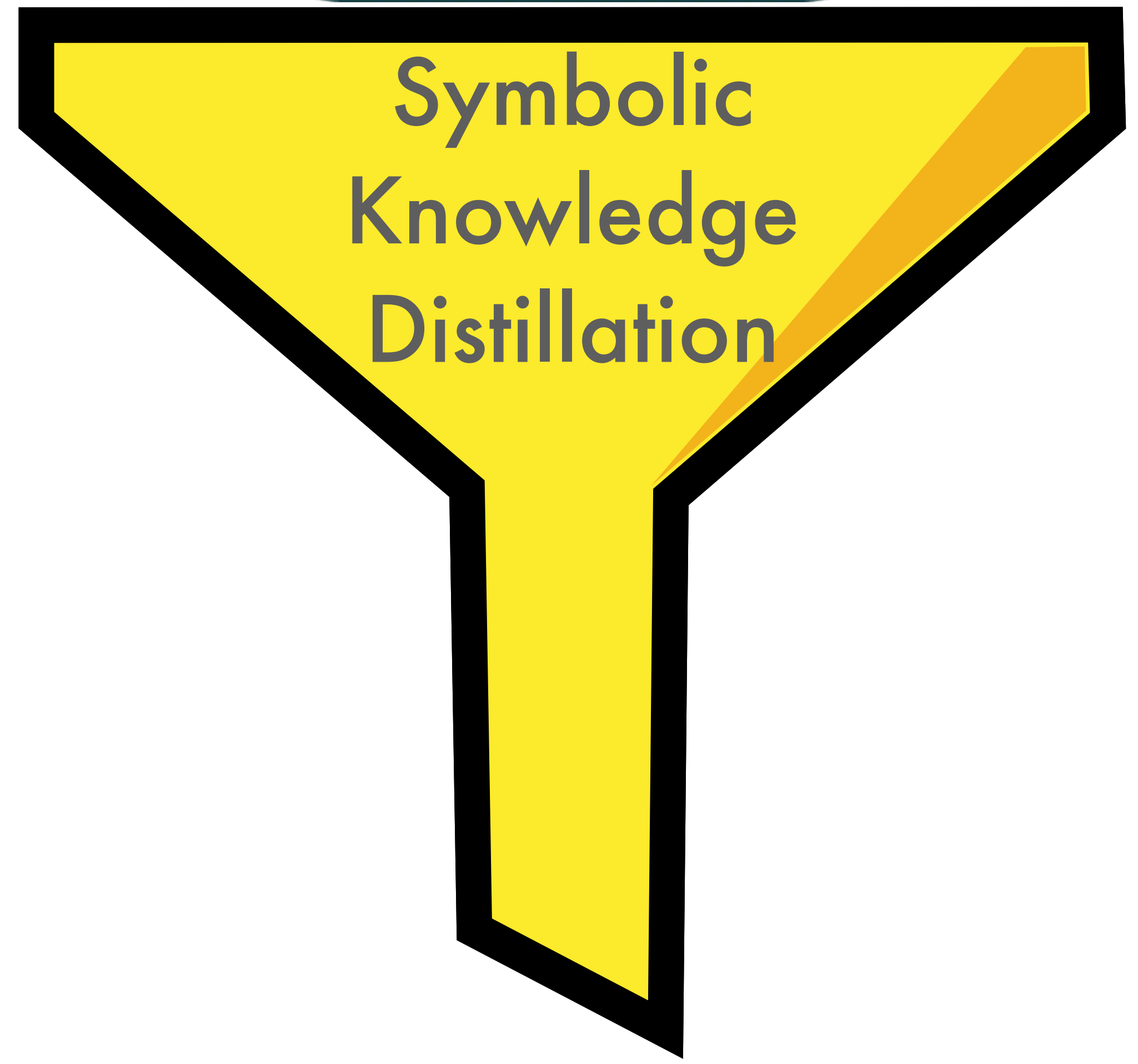




GPT-3



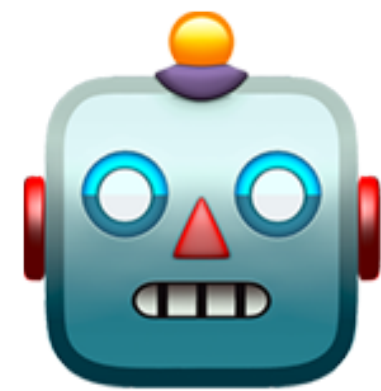




**SMALLER**

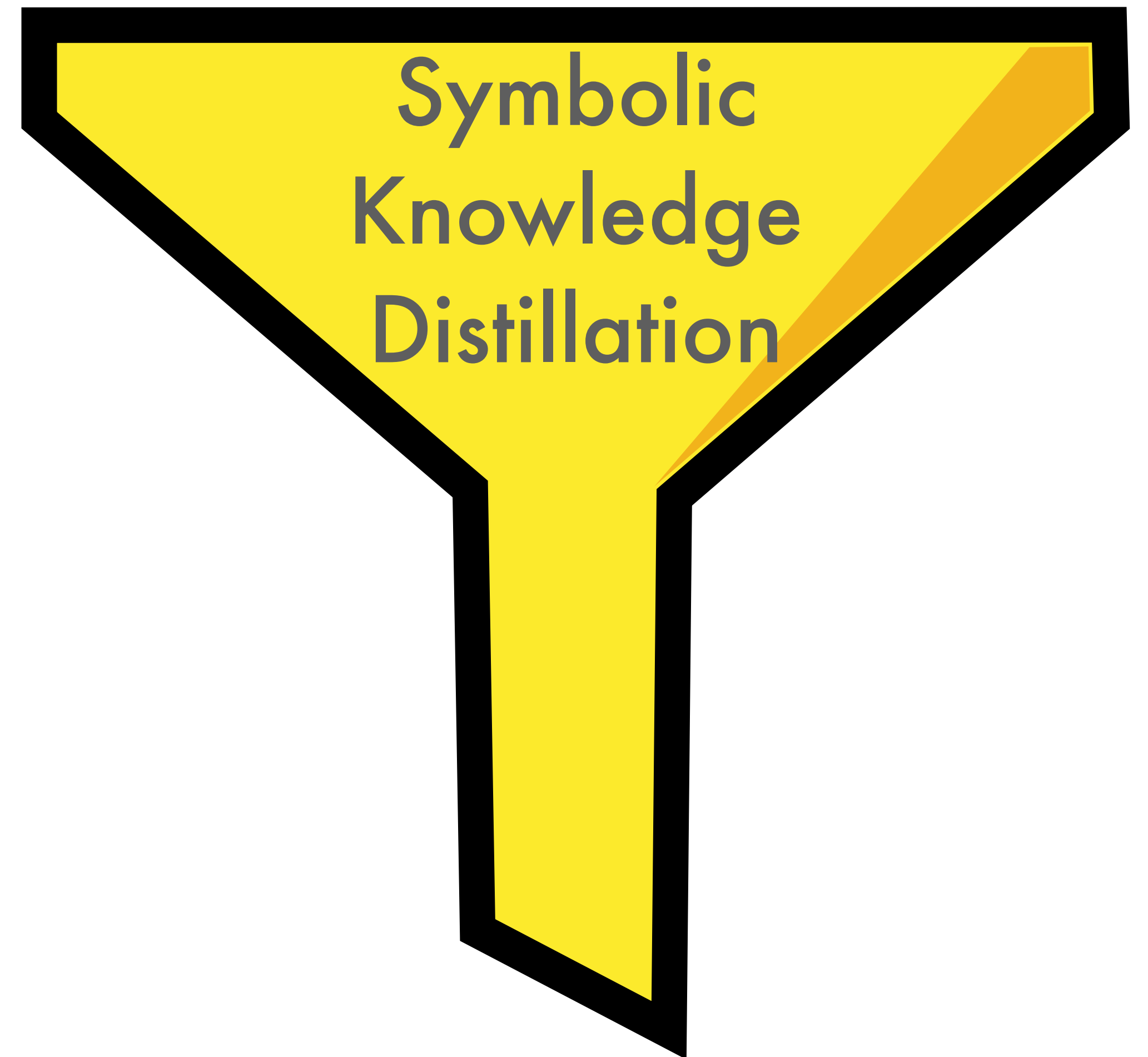
**AND**

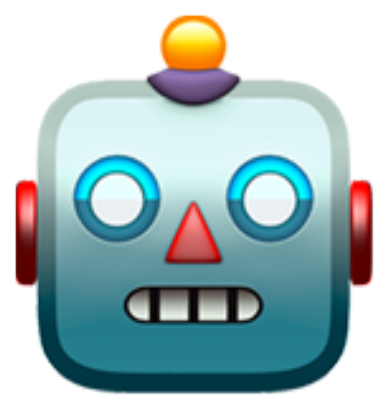
**BETTER**



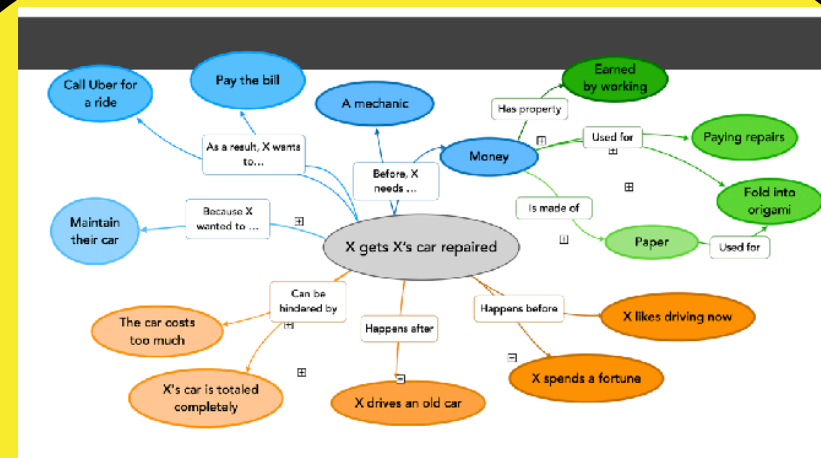
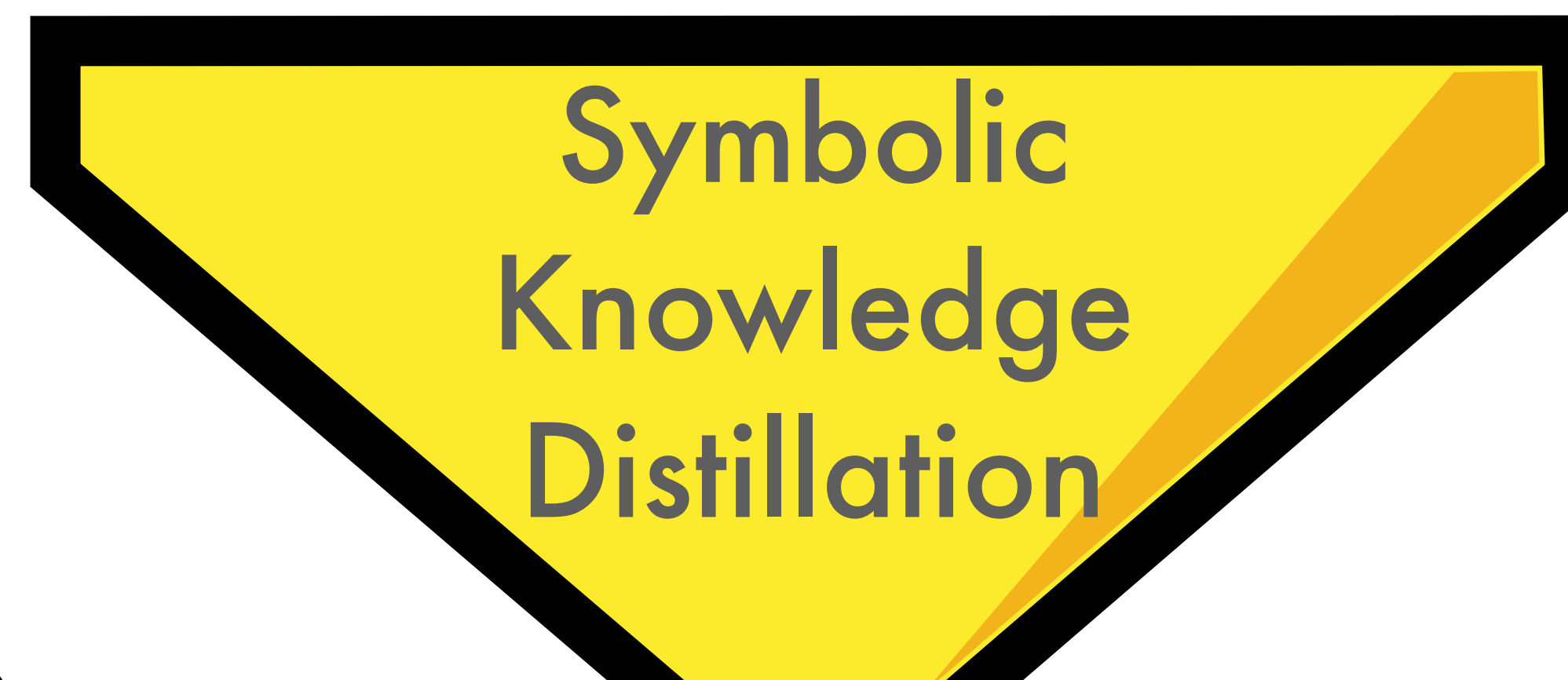
**EVEN POSSIBLE**

**???**





**Student Model**  
smaller & better



**Knowledge Graph**

6.5M high quality examples



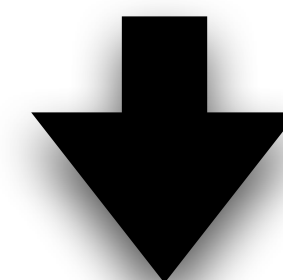
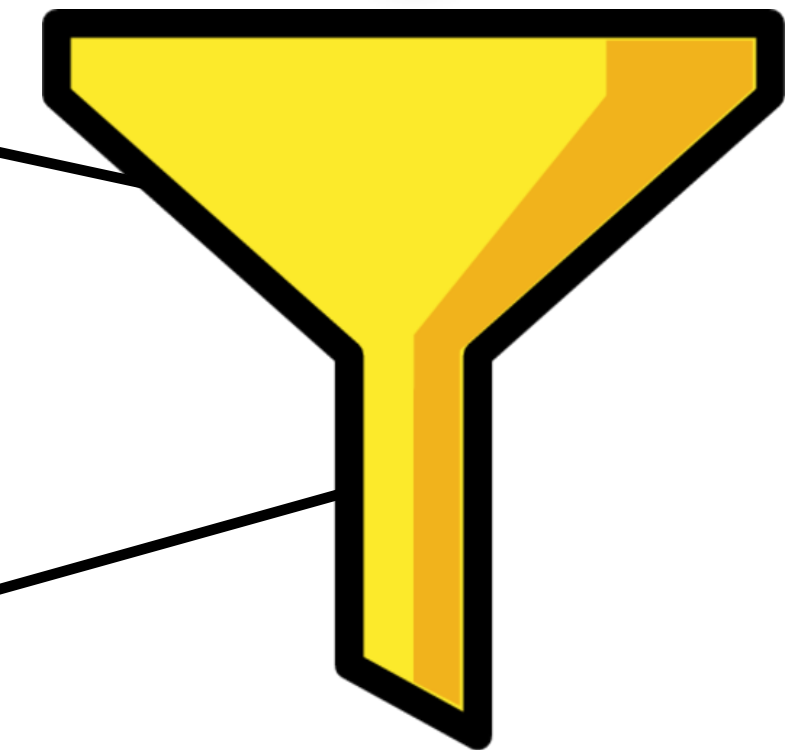
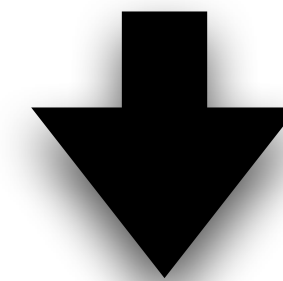
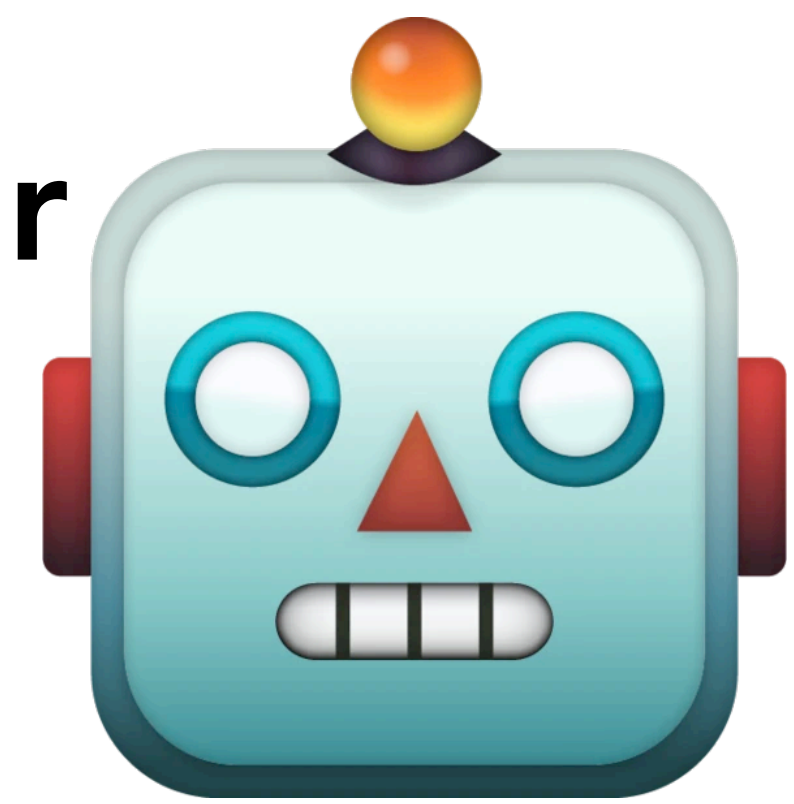
**Critic**

sort **good** and  
**bad** knowledge

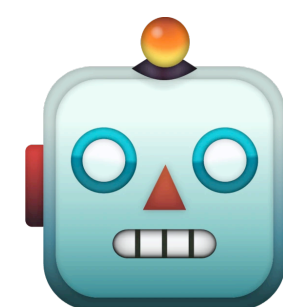
# Knowledge Distillation

(Hinton et al. 2015)

Teacher



Student

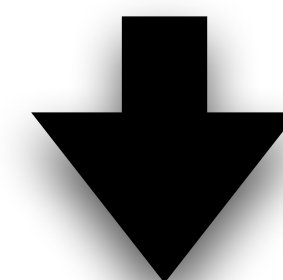
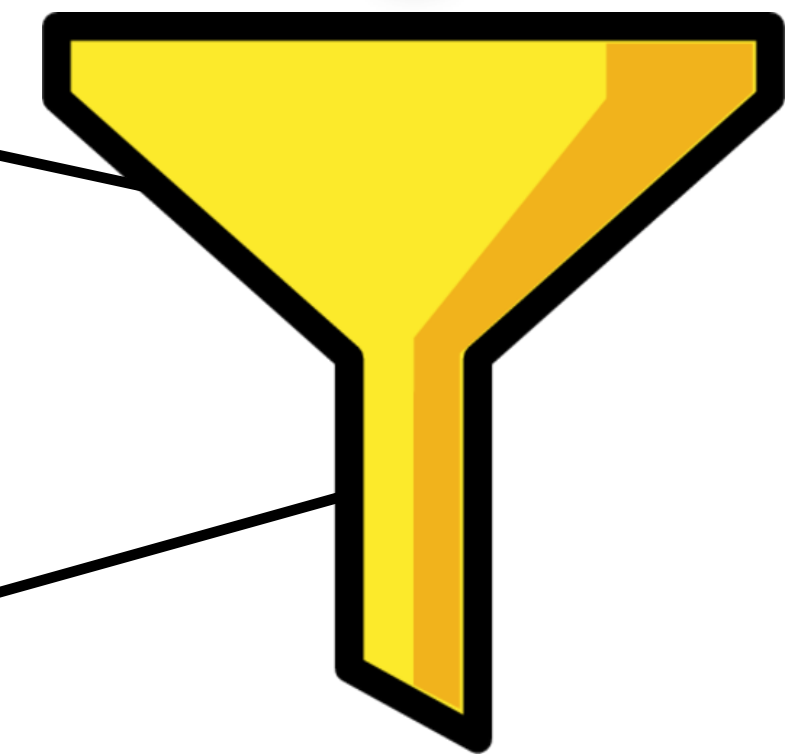
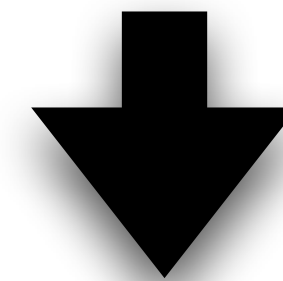
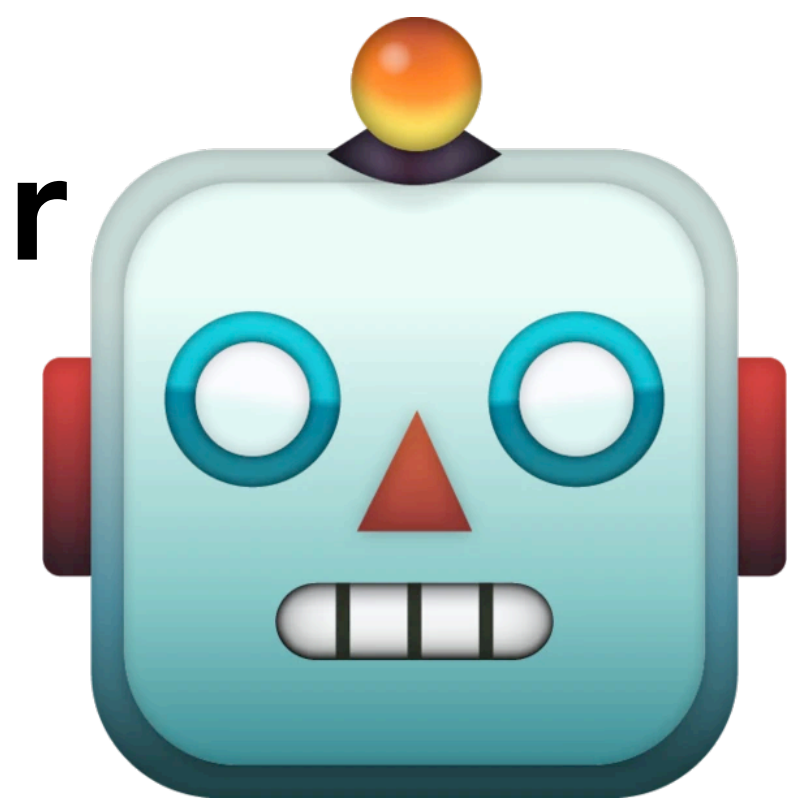


$$H(P_t, P_s) = - \sum_{y \in Y} P_t(y) \log P_s(y)$$

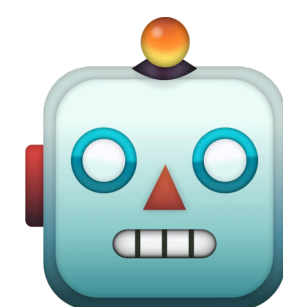
Train student to match  
teacher probabilities

# Symbolic Knowledge Distillation

Teacher



Student

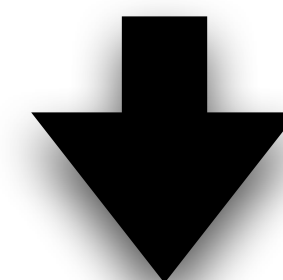
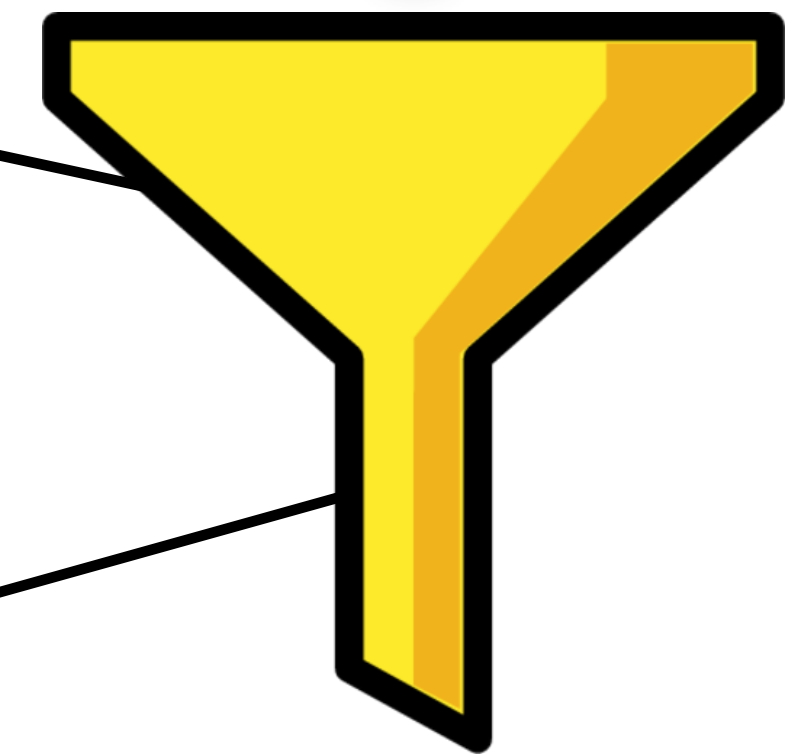
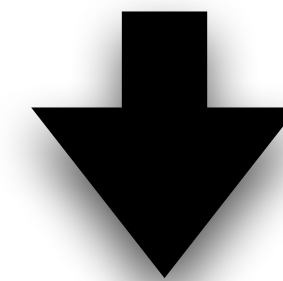
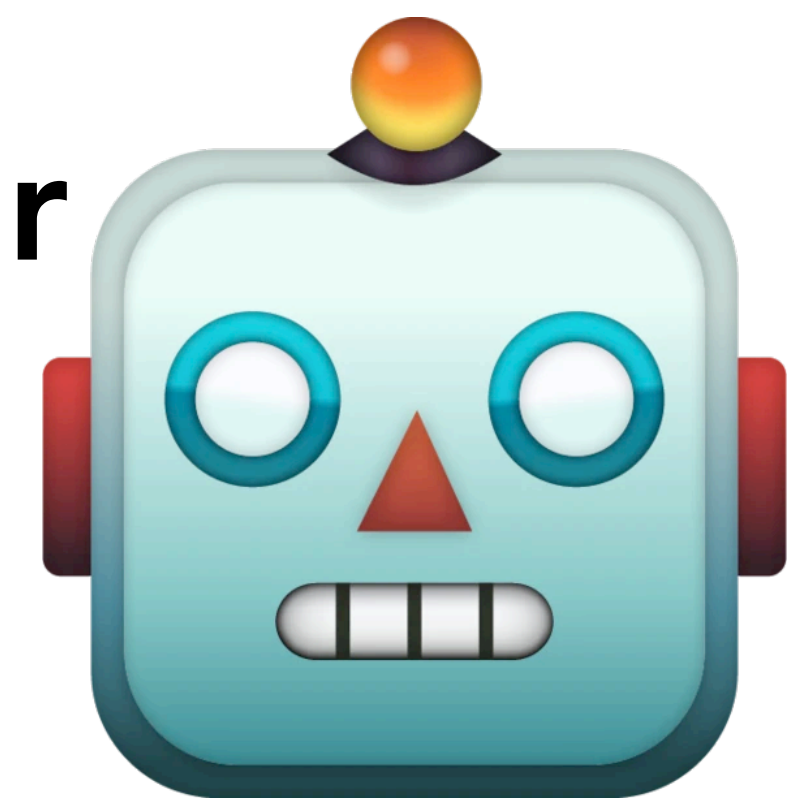


$$\cancel{H(P_t, P_s) = \sum_{y \in Y} P_t(y) \log P_s(y)}$$

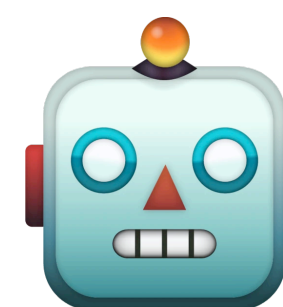
In generation,  $Y$  is all strings – intractable!

# Symbolic Knowledge Distillation

Teacher



Student



$$H(P_t, P_s) = \mathbb{E}_{y \sim P_t(y)} [-\log P_s(y)]$$

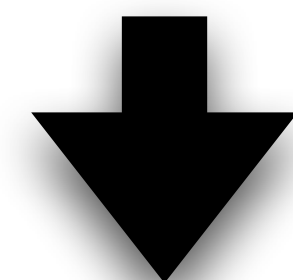
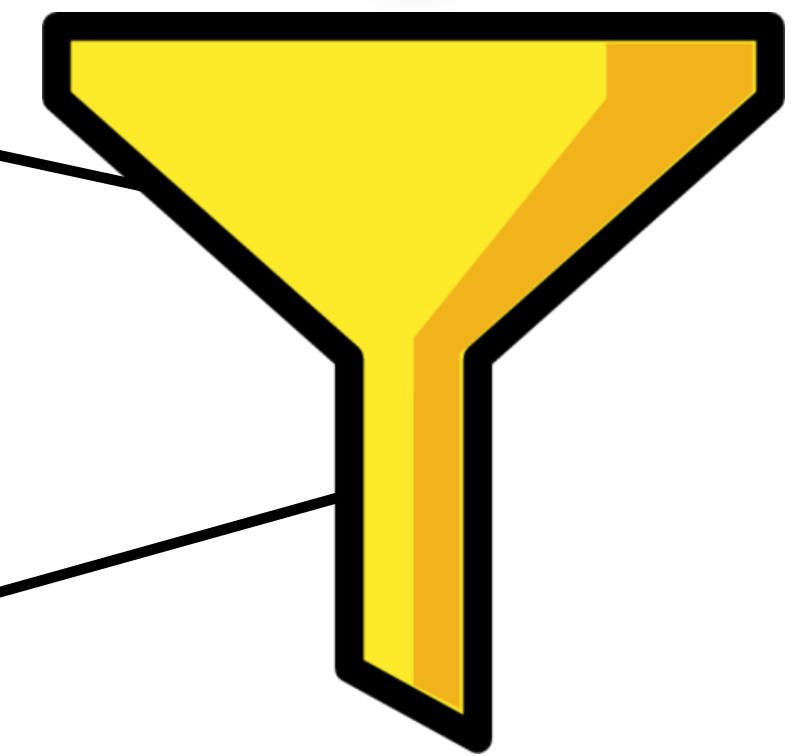
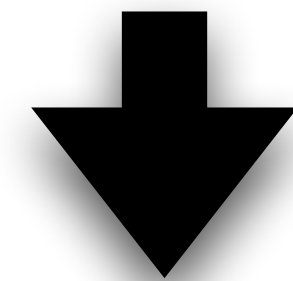
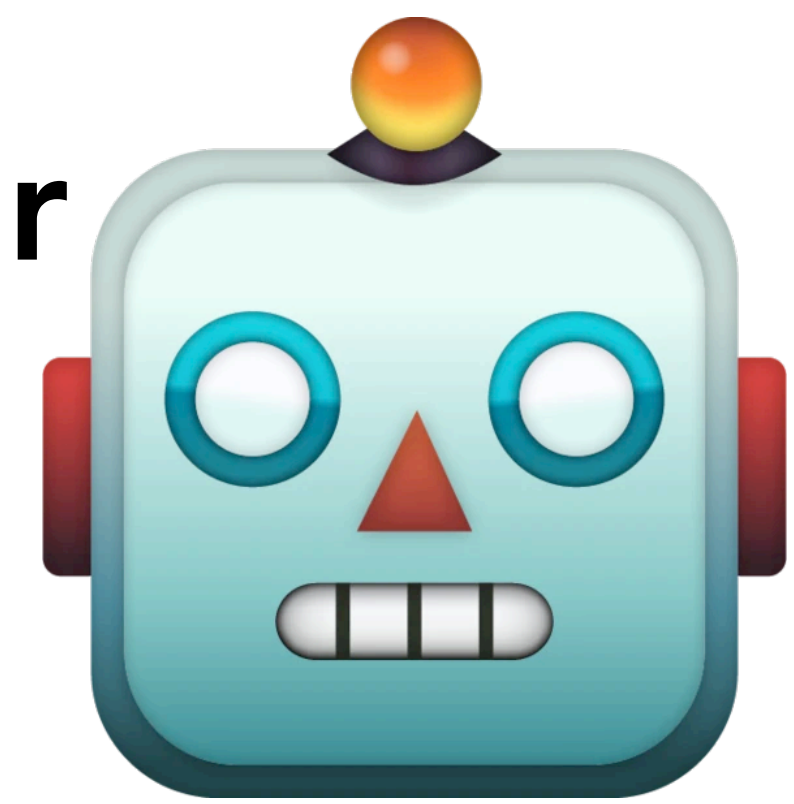
Estimate instead by  
generating examples!

Natural byproduct is a  
knowledge graph

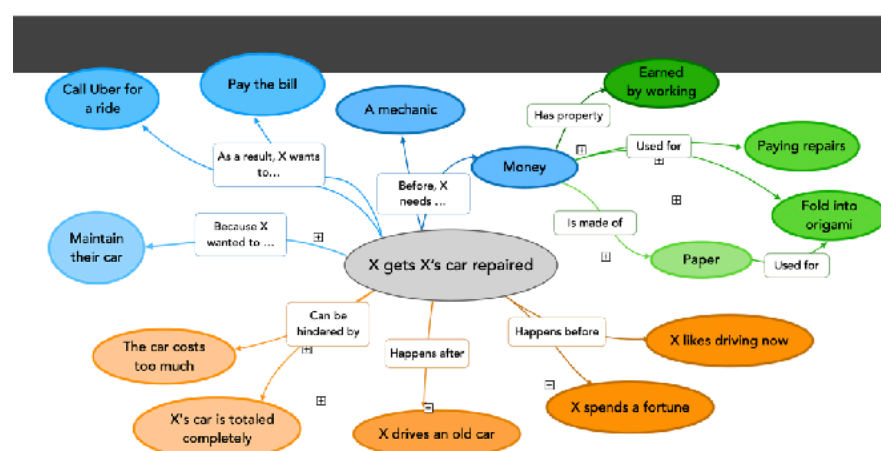


# Symbolic Knowledge Distillation

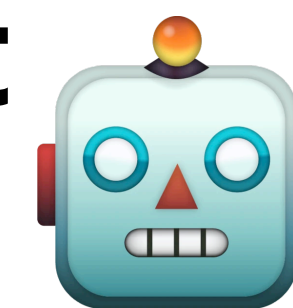
Teacher



KG



Student

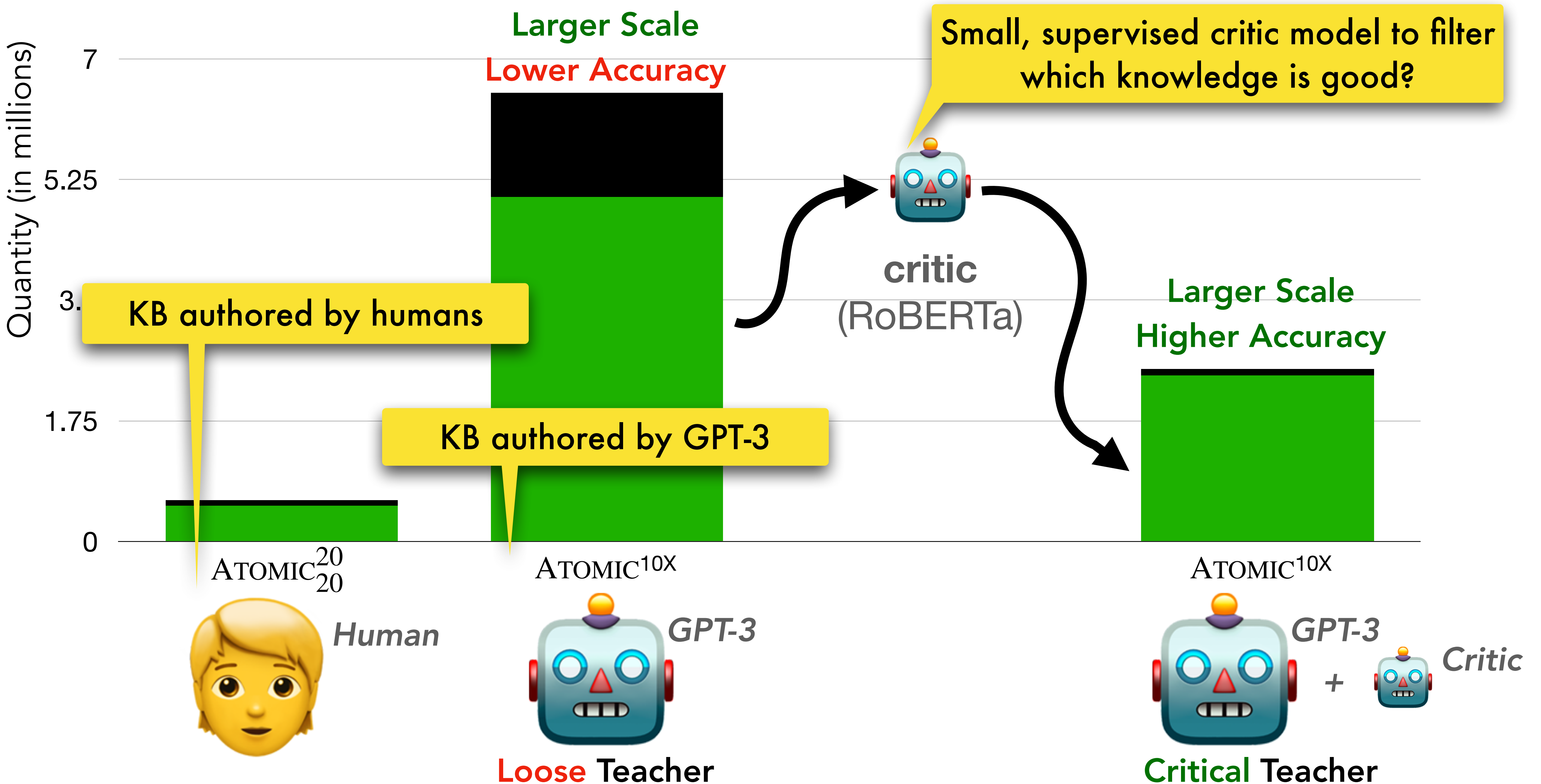


$$H(P_t, P_s) = \mathbb{E}_{y \sim P_t(y)} [-\log P_s(y)]$$

Estimate instead by generating examples!

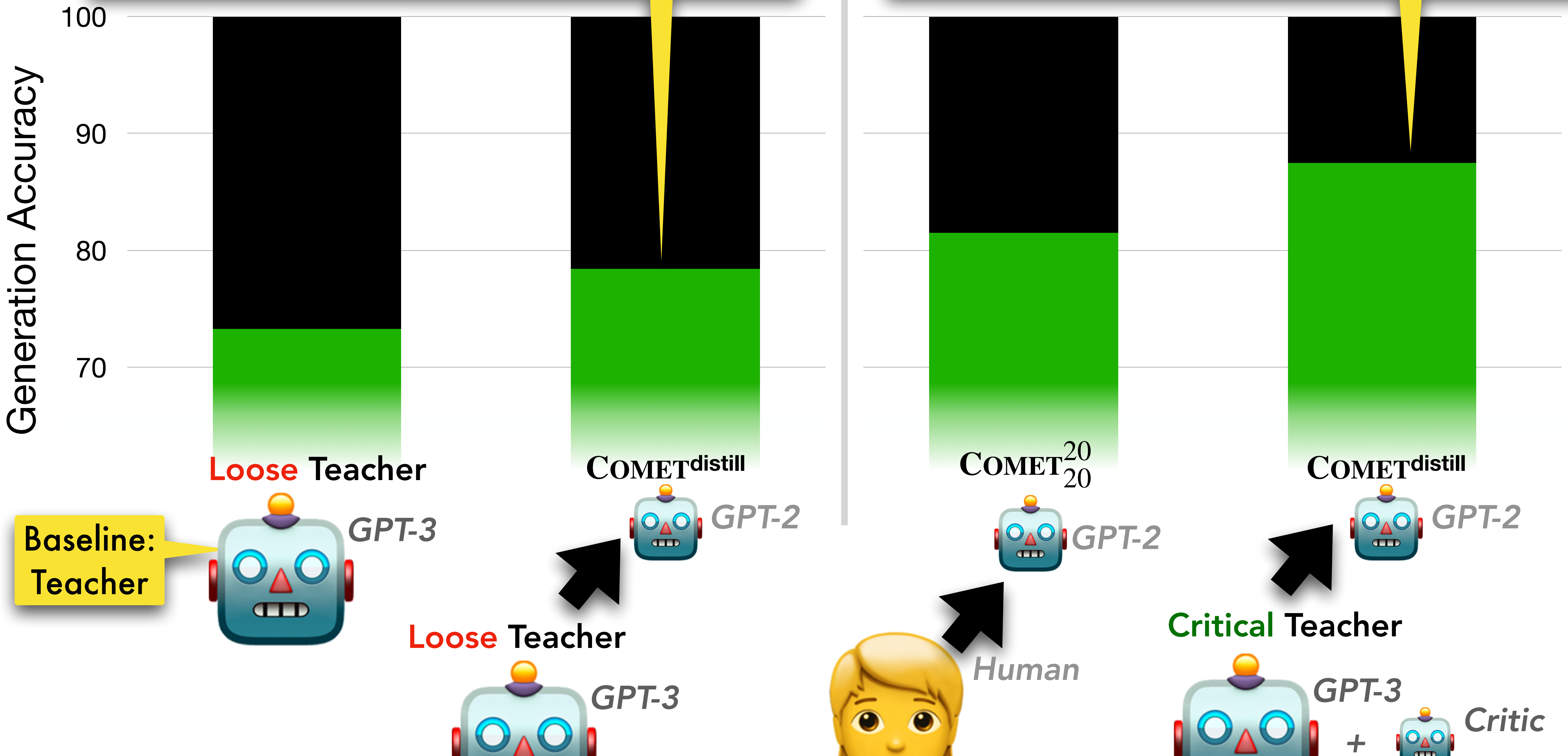
Natural byproduct is a knowledge graph

# Does Symbolic Knowledge Distillation Produce Good knowledge?



**Student COMET<sup>distill</sup> beats the teacher GPT-3 – smaller & better**

**Critical teacher results in a better student than human knowledge**

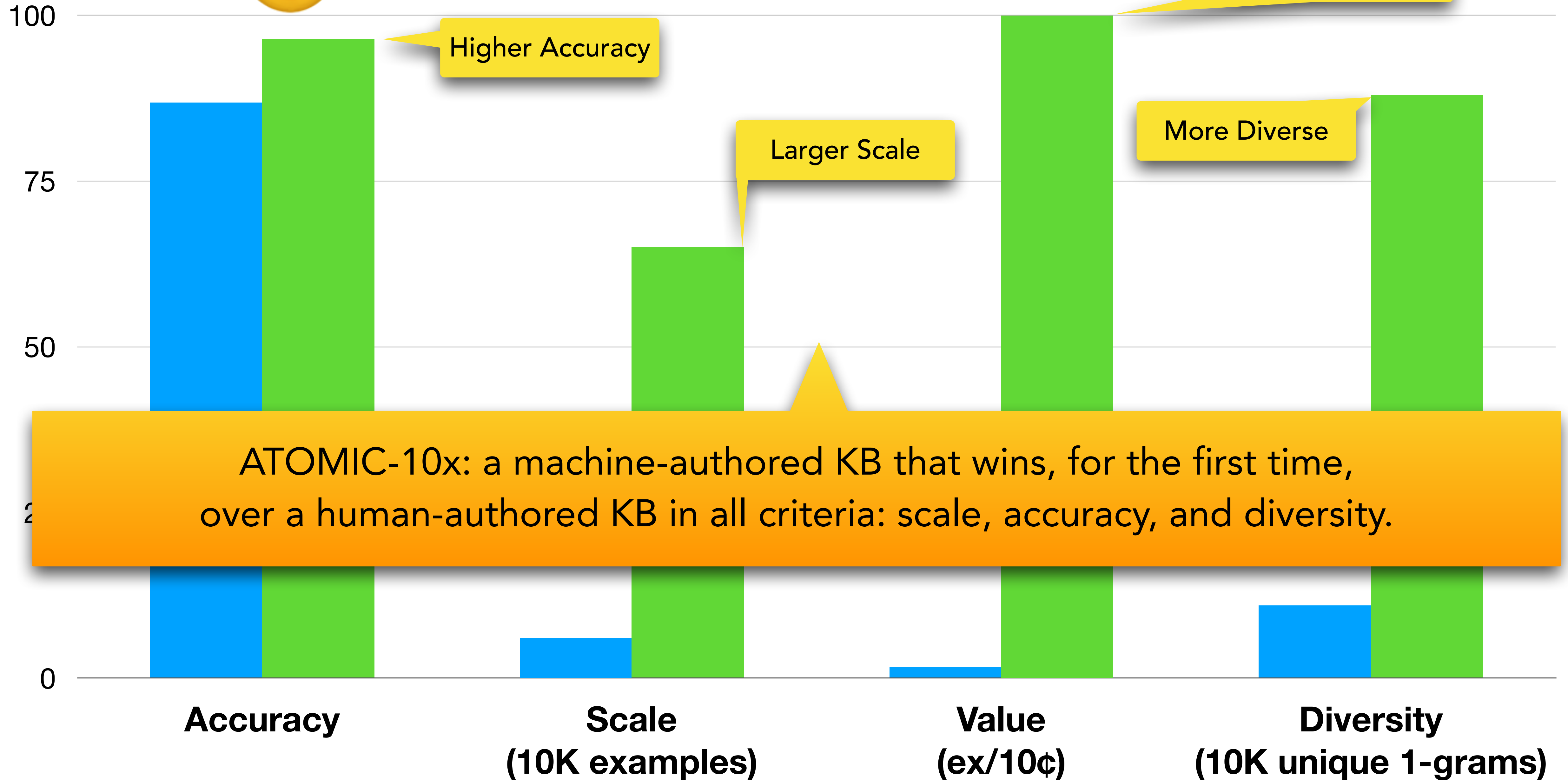




■ ATOMIC<sup>20</sup><sub>20</sub> VS



■ ATOMIC<sup>10X</sup>



Higher Accuracy

Larger Scale

More Diverse

Better Value

ATOMIC-10x: a machine-authored KB that wins, for the first time, over a human-authored KB in all criteria: scale, accuracy, and diversity.



Thanks! Questions?