

# Commonsense Reasoning: Models and New Challenges

Sean (Xiang) Ren

Department of Computer Science  
Information Science Institute  
USC

<http://inklab.usc.edu>



# Alibaba and Microsoft AI beat human scores on Stanford reading test

Neural networks edged past human scores on the measure of machine reading.

Bl

By

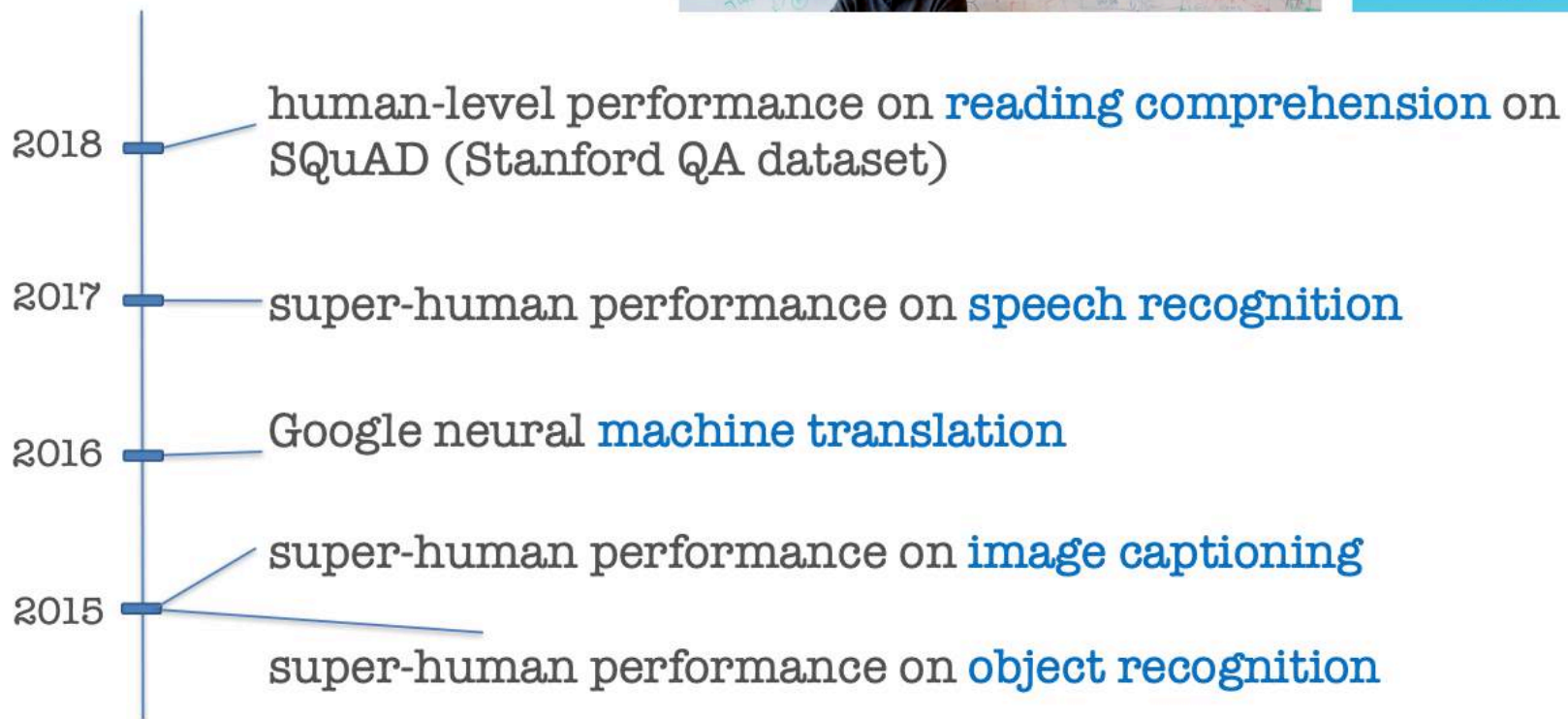
f



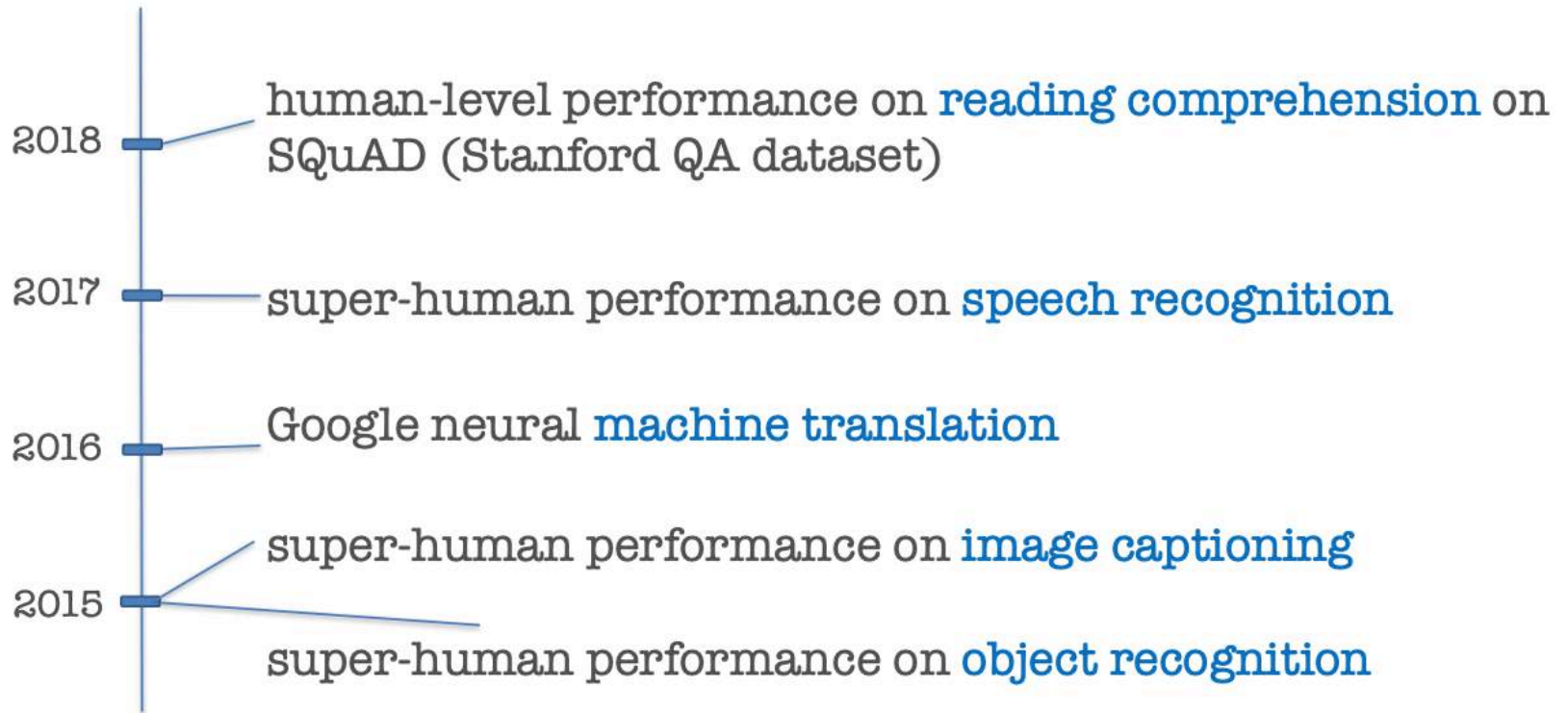
Rob LeFebvre, @roblef  
01.15.18 in [Personal Computing](#)

10  
Comments

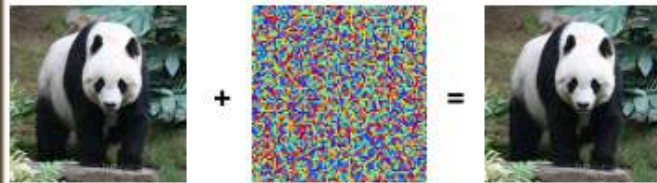
937  
Shares



# Done Solving AI?



# Solving a “dataset” vs the underlying “task”



Giant panda  
Object  
Recognition

Szegedy et al,  
2014....

Gibbon



VQA

Jabri et al,  
2017



A horse standing in the grass.

Captioning

MacLeod  
et al, 2017

How are you  
doing?



I don't know.

Dialogue

Li et al,  
2016



I don't know. I  
don't know. I  
don't know.

Open-ended

Generation

Holtzman  
et al, 2018

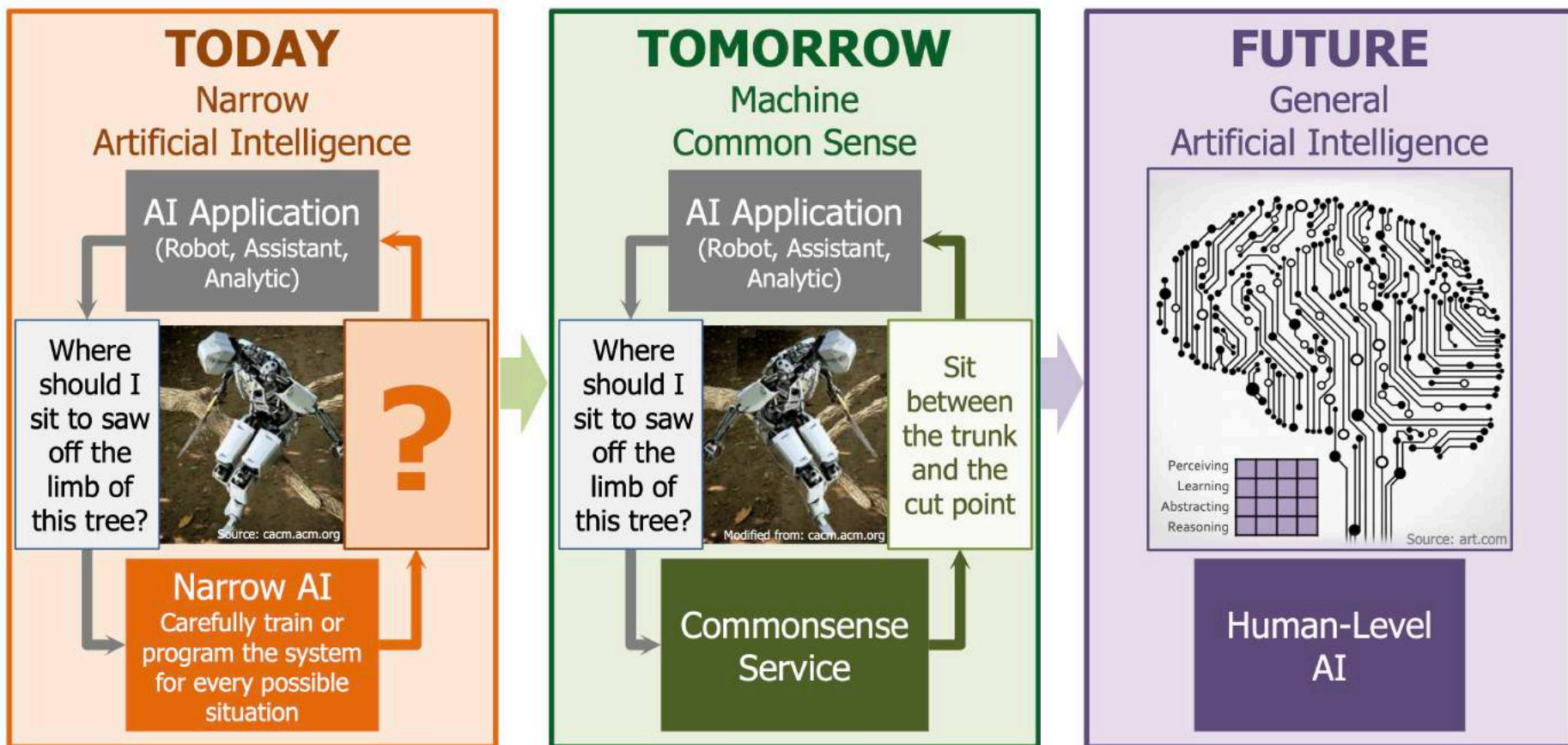
.... Nikola Tesla moved to  
Prague in 1880. ... **Tadakatsu**  
moved to **Chicago** in 1881.

Where did Tesla move in  
1880? **Chicago**

QA

Jia et al,  
2017

# Why Commonsense Knowledge?



# Commonsense problems in NLP

**NLU:** Multi-choice QA (w/o context)

Where do adults usually use glue sticks?

A: classroom    B: office    C: desk drawer

**NLG:** Constrained Sentence Generation (w/ a set of keywords)

Generate a daily-life scene about a concept-set: {apple, bag, tree}

*A boy picks some apples from a tree and puts them into a bag.*

# Commonsense Reasoning (CSR)?

- Definition of Common Sense: **the basic level of practical knowledge and reasoning**
  - Physical objects, properties, laws
  - Human behaviors / social conventions
  - Temporal commonsense
- The **human-like ability** to understand and generate everyday scenarios (situations, events)
- The **computation process** of manipulating commonsense knowledge to make compositional logical inference.

# This Talk

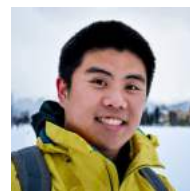
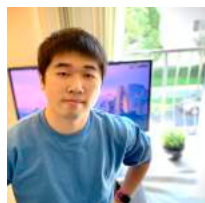
- Part I: **Discriminative** Commonsense Reasoning
  - Improving language understanding with commonsense
  - Models: **KagNet** and **multi-hop relational network**
- Part II: **Generative** Commonsense Reasoning
  - Imposing commonsense to text generation
  - A new task & dataset: **CommonGen**
  - Methods and Evaluation



## Part I

# KagNet: Knowledge-Aware Graph Networks for Commonsense Reasoning

*Bill Yuchen Lin*   *Xinyue Chen*   *Jamin Chen*   *Xiang Ren*



*University of Southern California - Information Science Institution*

INK Lab @ USC-ISI

<http://inklab.usc.edu>



USC



**EMNLP-IJCNLP 2019**  
Hong Kong, China



Association for  
Computational Linguistics

# Commonsense Question Answering

Where do adults usually use glue sticks?

A: classroom    B: office    C: desk drawer

What do you need to fill with ink to write notes on an A4 paper?

A: fountain pen    B: printer    C: pencil

Can you choose **the most plausible answer** based on daily life **commonsense** knowledge?

# Commonsense Question Answering

Where do adults usually use glue sticks?

A: classroom    B: office    C: desk drawer

What do you need to fill with ink to write notes on an A4 paper?

A: fountain pen    B: printer    C: pencil

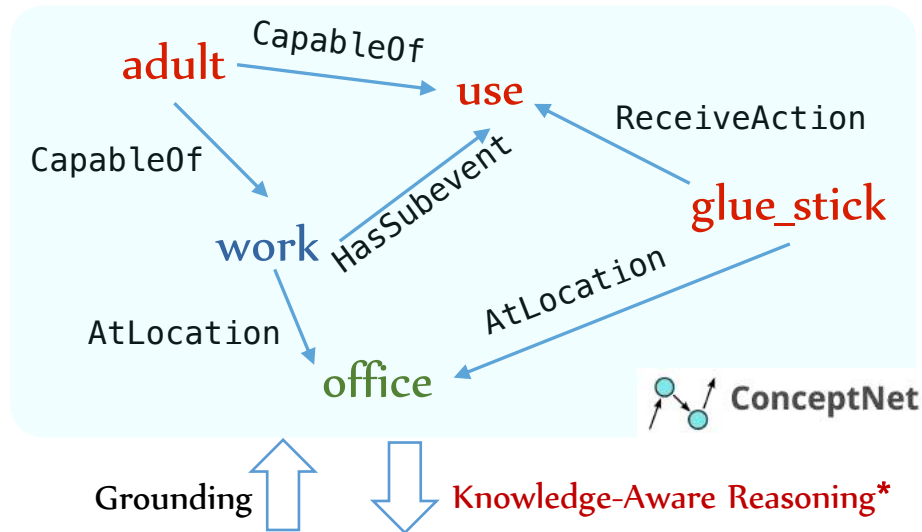
From the CommonsenseQA dataset (Talmor et al. NAACL 2019)

**Research question:**

**How can we impose commonsense in NLU models?**

# Knowledge-Aware Reasoning

Symbol Space



A Schema Graph for the choice B

Semantic Space

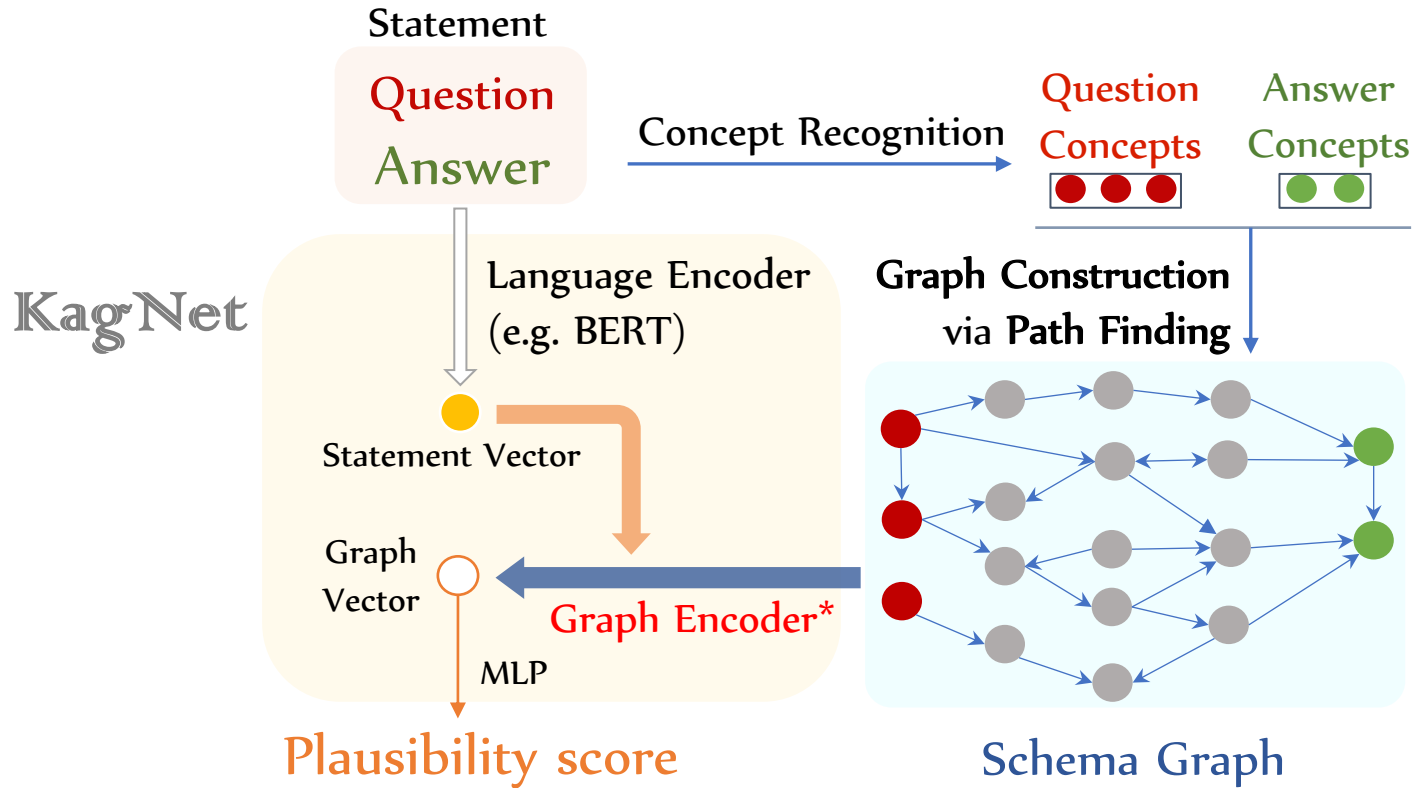
Where do adults use glue sticks?  
A: classroom B: office C: desk drawer

Question  
Answer Candidates

# Challenges in knowledge-aware reasoning

- **How can we find the schema graphs?**
  - Noisy and Incomplete
  - Numerous graphs; how to select the most related ones
- **How do we encode these graphs for reasoning?**
  - Complex multi-relational graph structures
  - **NO supervision in aligning** graphs and question-answer pairs
  - Need to be compatible with neural sentence encoders

# Proposed Framework Overview



# (1) Schema Graph Construction

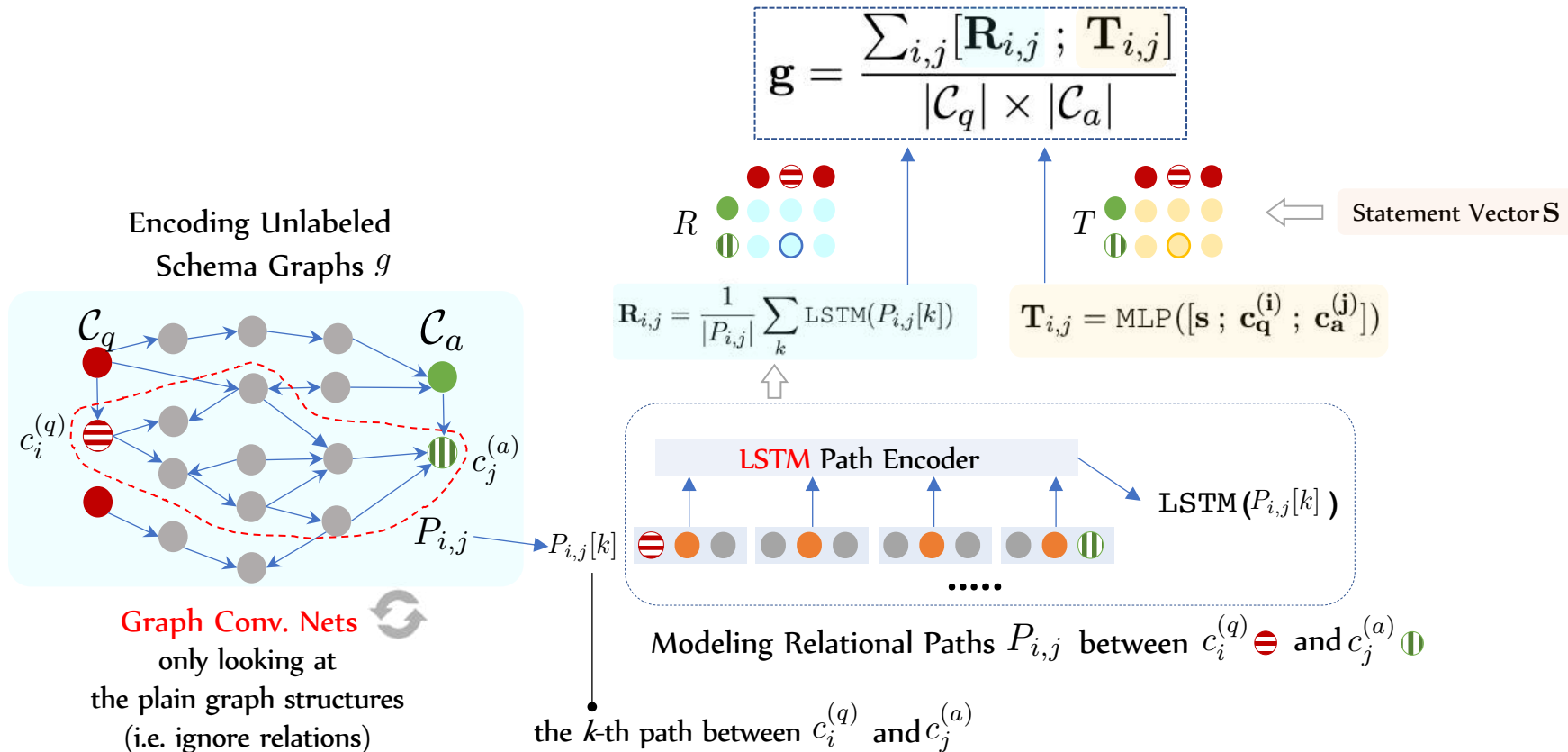
- **Concept Recognition**

- Tokenization / Lemmatization
- Match ConceptNet vocabulary
- Merge multiple smaller concepts into a longer one
  - e.g. " fountain", " pen" --> "fountain pen"
- **Question Concepts  $\mathcal{C}_q$**  and **Answer Concepts  $\mathcal{C}_a$**

- **Path Finding**

- Find paths between each QA-concept pair (one from  $\mathcal{C}_q$  and one from  $\mathcal{C}_a$ )
  - $\mathcal{P}_{i,j}$  denotes the set of paths between **i-th question concept** and **j-th answer concept**  
 $c_i^{(q)} \ominus , c_j^{(a)} \oplus$
- Path **pruning** by length ( $\leq 5$  nodes) and embedding-based metric.

# (2) Path-based Relational Graph Encoder





# (3) w/ Hierarchical Path-based Attention

- Two average pooling:

- Assuming all QA-concept pairs are equally important  $\mathbf{g} = \frac{\sum_{i,j} [\mathbf{R}_{i,j} ; \mathbf{T}_{i,j}]}{|\mathcal{C}_q| \times |\mathcal{C}_a|}$
- Assuming all paths are equally relevant  $\mathbf{R}_{i,j} = \frac{1}{|P_{i,j}|} \sum_k \text{LSTM}(P_{i,j}[k])$

- Modeling the two-level importance as latent weights:

$$\alpha_{(i,j,k)} = \mathbf{T}_{i,j} \mathbf{W}_1 \text{LSTM}(P_{i,j}[k]),$$

$$\hat{\alpha}_{(i,j,\cdot)} = \text{SoftMax}(\alpha_{(i,j,\cdot)}),$$

$$\hat{\mathbf{R}}_{i,j} = \sum_k \hat{\alpha}_{(i,j,k)} \cdot \text{LSTM}(P_{i,j}[k])$$

Path-Level Attention  
(attending on semantic space)

$$\beta_{(i,j)} = \mathbf{s} \mathbf{W}_2 \mathbf{T}_{i,j}$$

$$\hat{\beta}_{(\cdot,\cdot)} = \text{SoftMax}(\beta_{(\cdot,\cdot)})$$

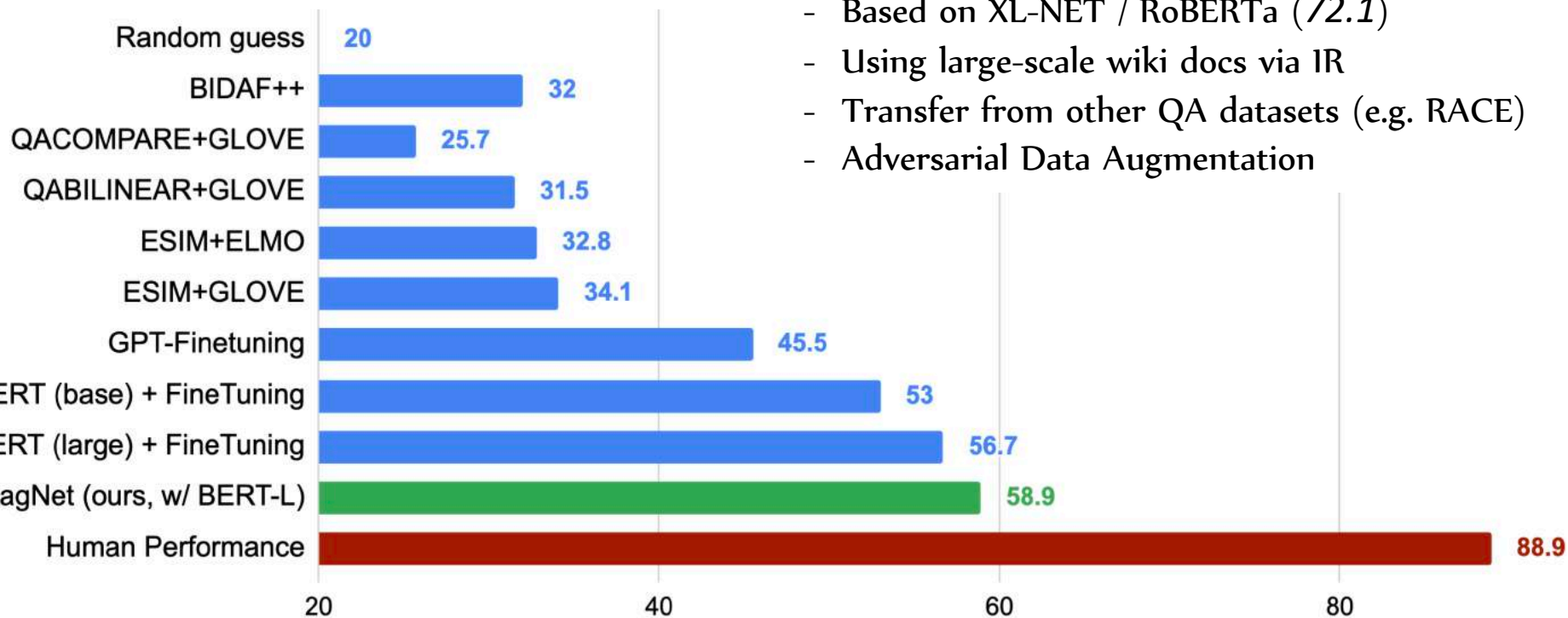
$$\hat{\mathbf{g}} = \sum_{i,j} \hat{\beta}_{(i,j)} [\hat{\mathbf{R}}_{i,j} ; \mathbf{T}_{i,j}]$$

ConceptPair-Level Attention  
(attending on statement)

# Experiments

Recent follow-up submissions:

- Based on XL-NET / RoBERTa (72.1)
- Using large-scale wiki docs via IR
- Transfer from other QA datasets (e.g. RACE)
- Adversarial Data Augmentation



More Performance on Official Test Set: <https://www.tau-nlp.org/csqa-leaderboard>

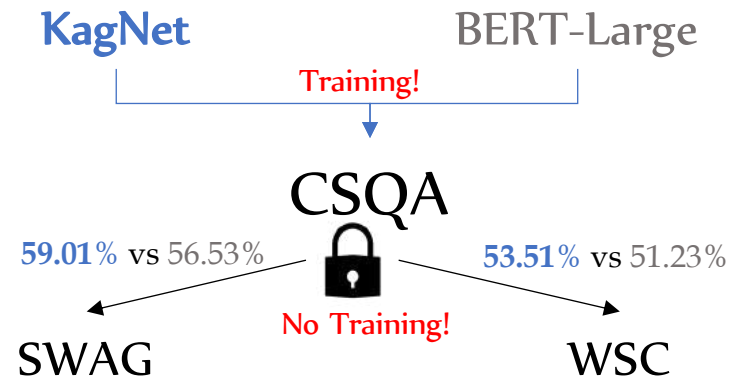
# Interpretability

What do you fill with ink to write on an A4 paper?

A: fountain pen ✓ (KagNet); B: printer (BERT);

C: squid D: pencil case (GPT); E: newspaper

# Transferability



# Conclusion

- A novel framework for knowledge-aware commonsense QA
- A graph neural network for **relational reasoning**.
  - GCN + Path-based LSTM + Hierarchical Attention
  - Promising for other reasoning tasks over graphs (e.g. GQA)
- Future directions in commonsense reasoning:
  - Towards **Learnable** Graph Construction (instead of heuristic algs.)
  - Explicitly deal with **negations** (“not”, “but”, etc.) and **comparisons** (“largest”, “most”, etc.).
    - Logical forms, executable semantic parsing.
  - **Interactively** reasoning over a sequence of questions
- Our code is at <https://github.com/INK-USC/KagNet>

# Multi-Hop Graph Relation Networks for Knowledge-Aware Question Answering

<https://arxiv.org/abs/2005.00646>

**Yanlin Feng<sup>♣\*</sup> Xinyue Chen<sup>♣\*</sup> Bill Yuchen Lin<sup>♥</sup> Peifeng Wang<sup>♥</sup> Jun Yan<sup>♥</sup> Xiang Ren<sup>♥</sup>**

fengyanlin@pku.edu.cn, kiwisher@sjtu.edu.cn,  
{yuchen.lin, peifengw, yanjun, xiangren}@usc.edu

<sup>♥</sup>University of Southern California

<sup>♣</sup>Peking University    <sup>♣</sup>Shanghai Jiao Tong University

# Motivation

- **KG-Augmented Commonsense QA:**

Leverage KG to provide knowledge which is not stated explicitly in the context.

1. Extract the paths/subgraph localized at the entities mentioned in the context from KG.
2. Encode the paths/subgraph.

- Previous works on encoding paths/sub-graph

  - **Path-based Modeling**

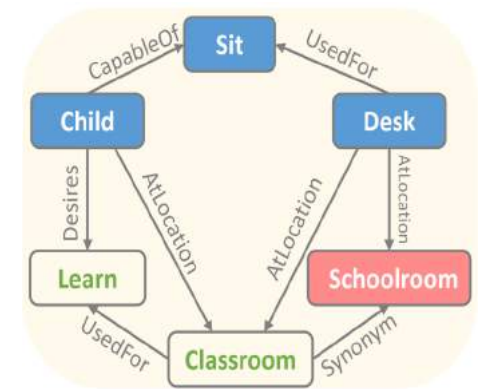
1. Model the relational paths with sequence model.
2. Use attention to aggregate the paths.

Interpretable, but not scalable.

  - **Relational Graph NN**

Model the subgraph with message passing.

Scalable, but lack transparency



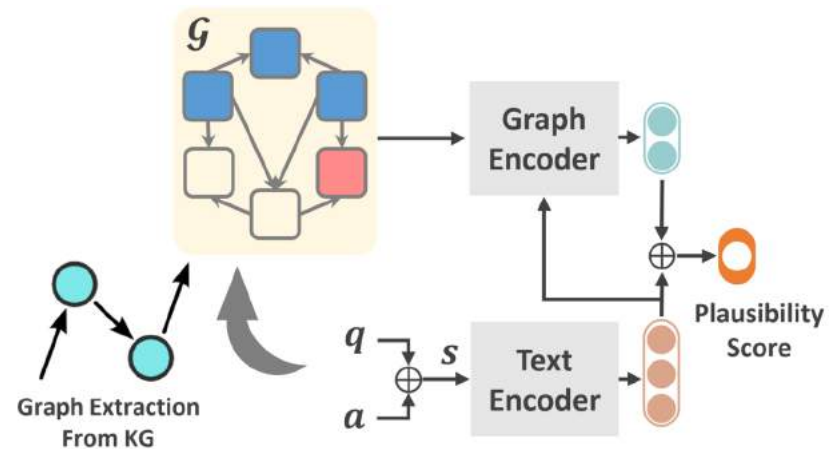
Where does a **child** likely **sit** at a **desk**?

- A. **Schoolroom**\*
- B. Furniture store
- C. Patio
- D. Office building
- E. Library

**Key idea: Modeling All Paths Directly in Graph Networks!**

# Reasoning Pipeline

- 1. Text Encoder:** Understand the textual input (question + answer choice).
- 2. Graph Encoder:** Reason over the contextual subgraphs.
- 3. Classifier:** Integrate the output from text/graph encoder to give a plausibility score.

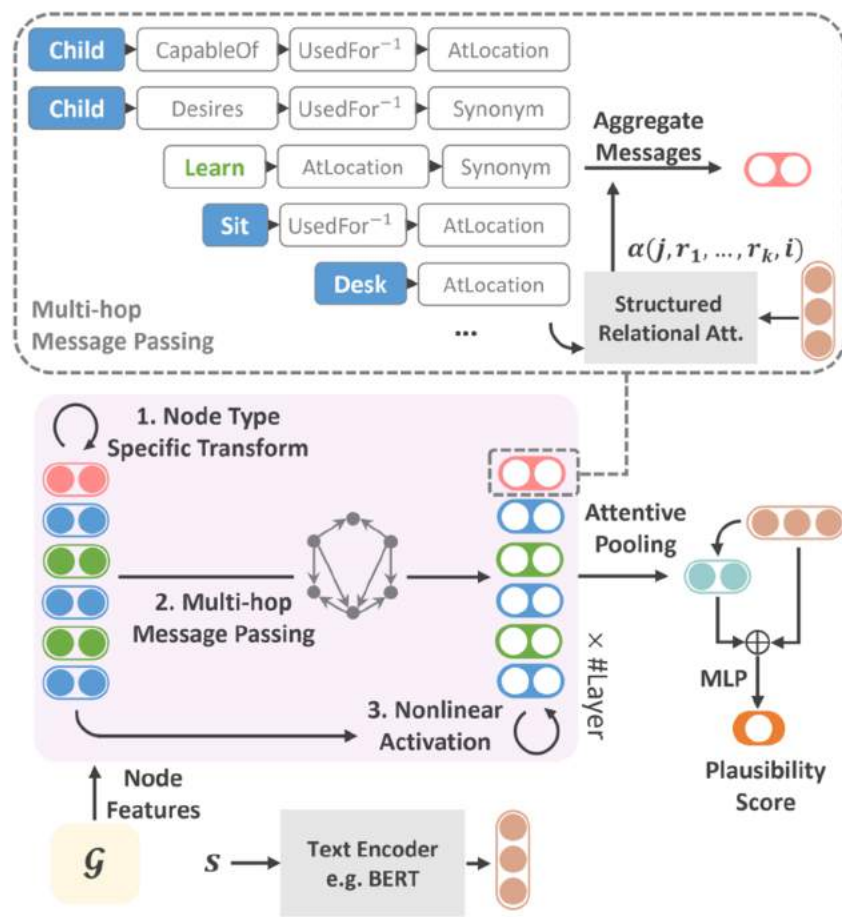


# Our Method for Encoding KG

**Goal:** To combine both interpretability (path-based modeling) and scalability (GNN).

**How:** Endow GNN with the capability to model paths **directly**.

1. Multi-Hop Message Passing
  - We extend message passing in GNN to k-hop paths modeling.
2. Structured Relational Attention
  - Incoming message for a node is aggregated by attention mechanism.





# Results

Methods	Single	Ensemble
RoBERTa <sup>†</sup>	72.1	72.5
RoBERTa + KEDGN <sup>†</sup>	72.5	74.4
RoBERTa + KE <sup>†</sup>	73.3	-
RoBERTa + HyKAS 2.0 <sup>†</sup> (Ma et al., 2019)	73.2	-
RoBERTa + FreeLB <sup>†</sup> (Zhu et al., 2020)	72.2	73.1
XLNet + DREAM <sup>†</sup>	66.9	73.3
XLNet + GR <sup>†</sup> (Lv et al., 2019)	75.3	-
ALBERT <sup>†</sup> (Lan et al., 2019)	-	<b>76.5</b>
RoBERTa + MHGRN ( $K = 2$ )	75.4	<b>76.5</b>

CommonsenseQA's Leaderboard

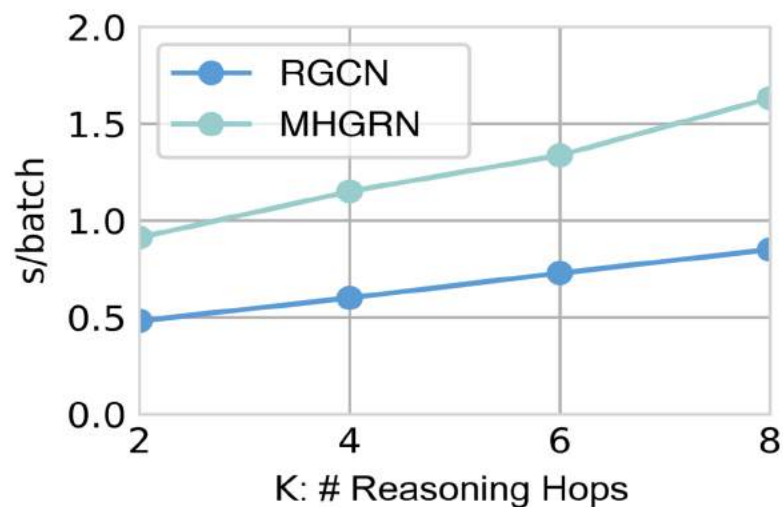
Methods	Dev (%)	Test (%)
T5-3B <sup>†</sup> (Raffel et al., 2019)	-	83.20
UnifiedQA <sup>†</sup> (Khashabi et al., 2020)	-	<b>87.20</b>
RoBERTa-Large (w/o KG)	66.76 ( $\pm 1.14$ )	64.80 ( $\pm 2.37$ )
+ RGCN	64.65 ( $\pm 1.96$ )	62.45 ( $\pm 1.57$ )
+ GconAttn	66.85 ( $\pm 1.82$ )	64.75 ( $\pm 1.48$ )
+ RN (1-hop)	64.85 ( $\pm 1.11$ )	63.65 ( $\pm 2.31$ )
+ RN (2-hop)	67.00 ( $\pm 0.71$ )	65.20 ( $\pm 1.18$ )
+ MHGRN ( $K = 3$ )	68.10 ( $\pm 1.02$ )	<b>66.85</b> ( $\pm 1.19$ )
AristoRoBERTaV7 <sup>†</sup>	79.2	77.8
+ MHGRN ( $K = 3$ )	78.6	<b>80.6</b>

OpenBookQA's Leaderboard

# Results

## Scalability

Model	Time	Space
<i><math>\mathcal{G}</math> is a dense graph</i>		
$K$ -hop KagNet	$\mathcal{O}(m^K n^{K+1} K)$	$\mathcal{O}(m^K n^{K+1} K)$
$K$ -layer RGCN	$\mathcal{O}(mn^2 K)$	$\mathcal{O}(mnK)$
MHGRN	$\mathcal{O}(m^2 n^2 K)$	$\mathcal{O}(mnK)$
<i><math>\mathcal{G}</math> is a sparse graph with maximum node degree <math>\Delta \ll n</math></i>		
$K$ -hop KagNet	$\mathcal{O}(m^K n K \Delta^K)$	$\mathcal{O}(m^K n K \Delta^K)$
$K$ -layer RGCN	$\mathcal{O}(mnK\Delta)$	$\mathcal{O}(mnK)$
MHGRN	$\mathcal{O}(m^2 n K \Delta)$	$\mathcal{O}(mnK)$

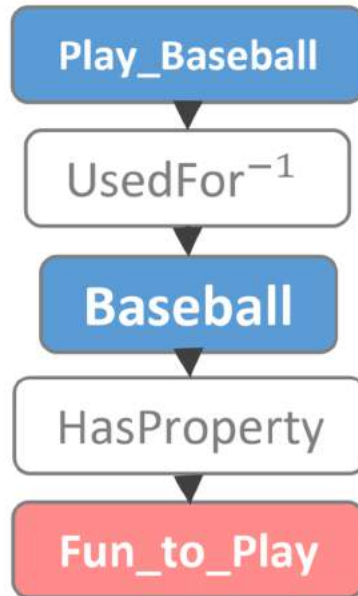


# Results

## Interpretability

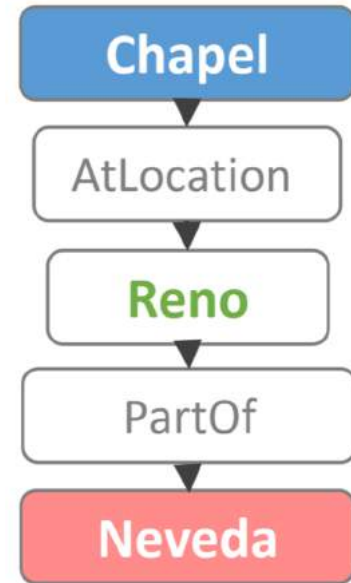
Why do parents encourage their kids to **play baseball**?

A. round B. cheap C. break window D. hard  
E. **fun to play**\*



Where is known for a multitude of wedding **chapels**?

A. town B. texas  
C. city  
D. church building  
E. **Nevada**\*



## Part II

# CommonGen: A Constrained Text Generation Challenge for Generative Commonsense Reasoning

<https://inklab.usc.edu/CommonGen/>

**Bill Yuchen Lin**♥   **Wangchunshu Zhou**♥   **Ming Shen**♥   **Pei Zhou**♥

**Chandra Bhagavatula**♠   **Yejin Choi**♠♠   **Xiang Ren**♥

♥ University of Southern California   ♠ Allen Institute for Artificial Intelligence

♠ Paul G. Allen School of Computer Science & Engineering, University of Washington



**USC** University of  
Southern California



**W**  
UNIVERSITY of  
WASHINGTON



# What is CommonGen?

- Most current tasks for machine commonsense focus on **discriminative** reasoning.
  - CommonsenseQA, SWAG.
- Humans not only use **commonsense knowledge** for understanding text, but also for **generating sentences**.

**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. 

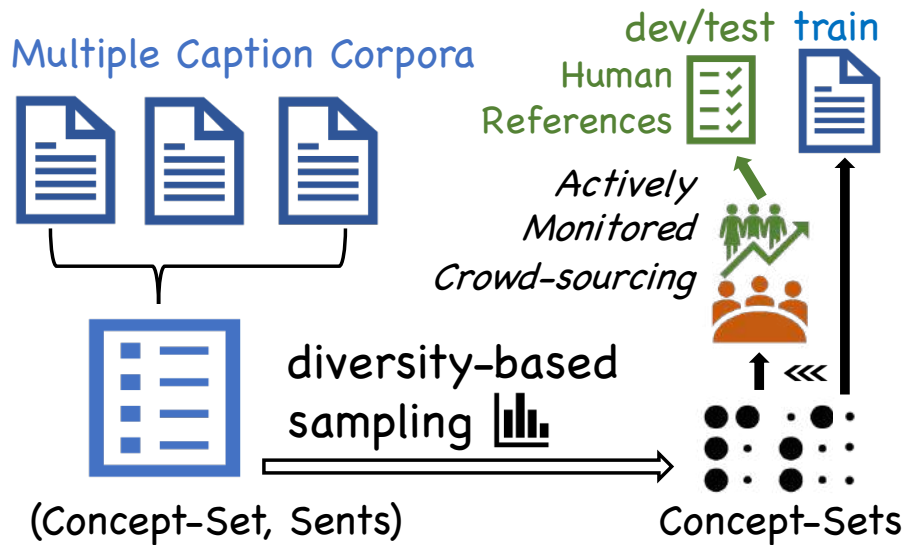
**Input:**

- A set of common concepts (actions & objects)

**Output:**

- A sentence that **describes an everyday scenario** the given concepts.

# Construction



Statistics	Train	Dev	Test
<b># Concept-Sets</b>	<b>32,651</b>	<b>993</b>	<b>1,497</b>
-Size = 3	25,020	493	-
-Size = 4	4,240	250	747
-Size = 5	3,391	250	750
<b># Sentences per Concept-Set Average Length</b>	67,389	4,018	6,042
	2.06	4.04	4.04
	10.54	11.55	13.34
<b># Unique Concepts</b>	4,697	766	1,248
<b># Unique Concept-Pairs</b>	59,125	3,926	8,777
<b># Unique Concept-Triples</b>	50,713	3,766	9,920
<b>% Unseen Concepts</b>	-	6.53%	8.97%
<b>% Unseen Concept-Pairs</b>	-	96.31%	100.00%
<b>% Unseen Concept-Triples</b>	-	99.60%	100.00%

# Why is it hard?

## Two key Challenges of CommonGen

(1) Relational knowledge are **latent** and **compositional**.

{ exercise, rope, wall, tie, wave }



### Underlying Relational Commonsense Knowledge

(exercise, HasSubEvent , releasing energy)

(rope, UsedFor, tying something)

(releasing energy, HasPrerequisite, motion)

(wave, IsA, motion) ; (rope, UsedFor, waving)

The motion costs more energy if ropes are tied to a wall.



### Relational Reasoning for Generation

A woman in a gym exercises by waving ropes tied to a wall.

Category	Relations	1-hop	2-hop
<i>Spatial knowledge</i>	AtLocation, LocatedNear	9.40%	39.31%
<i>Object properties</i>	UsedFor, CapableOf, PartOf, ReceivesAction, MadeOf, FormOf, HasProperty, HasA	9.60%	44.04%
<i>Human behaviors</i>	CausesDesire, MotivatedBy, Desires, NotDesires, Manner	4.60%	19.59%
<i>Temporal knowledge</i>	Subevent, Prerequisite, First/Last-Subevent	1.50%	24.03%
<i>General</i>	RelatedTo, Synonym, DistinctFrom, IsA, HasContext, SimilarTo	74.89%	69.65%

# Why is it hard?

## Two key Challenges of CommonGen

(2) Compositional Generalization for unseen concept compounds.

	Training
$x_1 = \{ \text{apple, bag, put} \}$	
$y_1 = \text{a girl puts an apple in her bag}$	
$x_2 = \{ \text{apple, tree, pick} \}$	
$y_2 = \text{a man picks some apples from a tree}$	
$x_3 = \{ \text{apple, basket, wash} \}$	
$y_3 = \text{a boy takes an apple from a basket and washes it.}$	

↓ Compositional Generalization

$x = \{ \text{pear, basket, pick, put, tree} \}, y = ?$	
<b>Reference:</b> "a girl picks some pear from a tree and put them in her basket."	Test

▶ Unseen Concept in Training

Statistics	Train	Dev	Test
<b># Concept-Sets</b>	<b>32,651</b>	<b>993</b>	<b>1,497</b>
-Size = 3	25,020	493	-
-Size = 4	4,240	250	747
-Size = 5	3,391	250	750
<b># Sentences per Concept-Set</b>	67,389	4,018	6,042
<b>Average Length</b>	2.06	4.04	4.04
	10.54	11.55	13.34
<b># Unique Concepts</b>	4,697	766	1,248
<b># Unique Concept-Pairs</b>	59,125	3,926	8,777
<b># Unique Concept-Triples</b>	50,713	3,766	9,920
<b>% Unseen Concepts</b>	-	6.53%	8.97%
<b>% Unseen Concept-Pairs</b>	-	96.31%	100.00%
<b>% Unseen Concept-Triples</b>	-	99.60%	100.00%



# Experimental Results

Model \ Metrics	ROUGE-2/L		BLEU-3/4		METEOR	CIDEr	SPICE	Coverage	
bRNN-CopyNet (Gu et al., 2016)	7.61	27.79	10.70	5.70	15.80	4.79	15.00	51.15	(1) Seq2seq models
Trans-CopyNet	8.78	28.08	11.90	7.10	15.50	4.61	14.60	49.06	
MeanPooling-CopyNet	9.66	31.14	10.70	6.10	16.40	5.06	17.20	55.70	
LevenTrans. (Gu et al., 2019)	10.58	32.23	19.70	11.60	20.10	7.54	19.00	63.81	
ConstLeven. (Susanto et al., 2020)	11.82	33.04	18.90	10.10	24.20	10.51	22.20	94.51	
GPT-2 (Radford et al., 2019)	17.18	39.28	30.70	21.10	26.20	12.15	25.90	79.09	(2) Fine-tuning pre-trained LMs
BERT-Gen (Bao et al., 2020)	18.05	40.49	30.40	21.10	27.30	12.49	27.30	86.06	
UniLM (Dong et al., 2019)	21.48	<b>43.87</b>	<u>38.30</u>	<u>27.70</u>	29.70	<u>14.85</u>	30.20	89.19	
UniLM-v2 (Bao et al., 2020)	18.24	40.62	31.30	22.10	28.10	13.10	28.10	89.13	
BART (Lewis et al., 2019)	<b>22.23</b>	41.98	36.30	26.30	<b>30.90</b>	13.92	<u>30.60</u>	<b>97.35</b>	
T5-Base (Raffel et al., 2019)	14.57	34.55	26.00	16.40	23.00	9.16	22.00	76.67	
T5-Large (Raffel et al., 2019)	<u>22.01</u>	<u>42.97</u>	<b>39.00</b>	<b>28.60</b>	<u>30.10</u>	<b>14.96</b>	<b>31.60</b>	<u>95.29</u>	
Human Performance	48.88	63.79	48.20	44.90	36.20	43.53	63.50	99.31	(3) Agreement

Manual Eval. →

	C.Leven	GPT	BERT-G.	UniLM	BART	T5
Hit@1	3.2	21.5	22.3	21.0	<u>26.3</u>	<b>26.8</b>
Hit@3	18.2	63.0	59.5	<u>69.0</u>	<u>69.0</u>	<b>70.3</b>
Hit@5	51.4	95.5	95.3	<u>96.8</u>	<u>96.3</u>	<b>97.8</b>

# Case Study & Transfer Learning

Concept-Set: { hand, sink, wash, soap }

**[bRNN-CopyNet]:** a hand works in the sink .

**[MeanPooling-CopyNet]:** the hand of a sink being washed up

**[ConstLeven]:** a hand strikes a sink to wash from his soap.

**[GPT-2]:** hands washing soap on the sink.

**[BERT-Gen]:** a woman washes her hands with a sink of soaps.

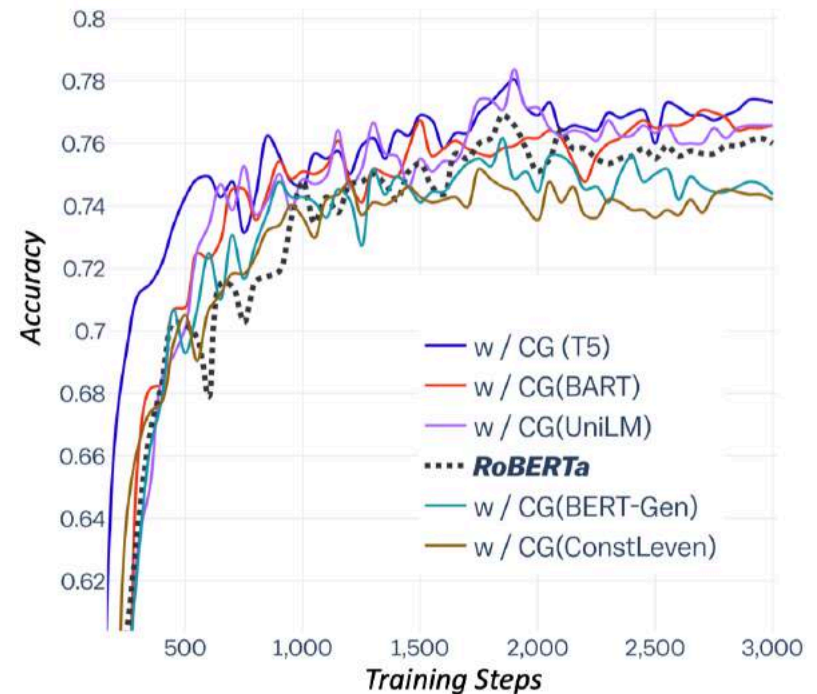
**[UniLM]:** hands washing soap in the sink

**[BART]:** a man is washing his hands in a sink with soap and washing them with hand soap.

**[T5]:** hand washed with soap in a sink.



1. A girl is washing her hands with soap in the bathroom sink.
2. I will wash each hand thoroughly with soap while at the sink.
3. The child washed his hands in the sink with soap.
4. A woman washes her hands with hand soap in a sink.
5. The girl uses soap to wash her hands at the sink.



**Learning curve for the transferring study** (acc on dev). We use trained CommonGen models to generate choice-specific context for the CommonsenseQA task.

# Learning with Natural Language Explanations

Sentiment on ENT is  
**positive** or **negative**?

*Users' natural language  
explanations*

$x_1$ : There was a long wait for a table outside, but it was a little too hot in the sun anyway so our ENT was very nice.



**Positive**, because the words “*very nice*” is within 3 words after the ENT.

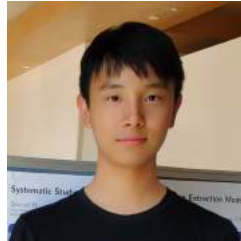
**Relation** between ENT1 and ENT2?

$x_2$ : Officials in Mumbai said that the two suspects, David Headley, and ENT1, who was born in Pakistan but is a ENT2 citizen, both visited Mumbai before the attacks.



**per: nationality**, because the words “*is a*” appear right before ENT2 and the word “*citizen*” is right after ENT2.

## Students



## Research Partnership



## Collaborators

Dan MacFarland, Sociology, Stanford University  
Jure Leskovec, Computer Science, Stanford University  
Dan Jurafsky, Computer Science, Stanford University  
Jiawei Han, Computer Science, UIUC  
Morteza Dehghani, Psychology, USC  
Kenneth Yates, Clinical Education, USC  
Craig Knoblock, USC ISI  
Curt Langlotz, Bioinformatics, Stanford University  
Kuansan Wang, Microsoft Academic  
Leonardo Neves, Snap Research  
Mark Musen, Bioinformatics, Stanford University

## Funding



IARPA  
BE THE FUTURE

J.P.Morgan



Google

SCHMIDT FAMILY  
FOUNDATION

amazon



Adobe

# Thank you!

- USC Intelligence and Knowledge Discovery (INK) Lab
  - <http://inklab.usc.edu/>
- Code: <https://github.com/INK-USC>
  - [xiangren@usc.edu](mailto:xiangren@usc.edu)
  - @xiangrenNLP

