Feb 2020

Nasrin Mostafazadeh

S@nasrinmmm

CONTRACTOR CONTRACTOR

Towards AI systems that can build coherent causal models of what they read!

State of Artificial Intelligence, ~15 years ago

RoboCup Competitions



https://www.youtube.com/watch?v=YPYVL5FpS6s

Classic Motivating NLU Problem

Deemed very challenging for AI systems at the time!

- The monkey ate the banana because it was hungry.
 - Question: What is it? Monkey or the banana?
 - Correct answer: Monkey

State of Artificial Intelligence, NOW!

Boston Dynamics' Robots

(2019)



Stanford CoreNLP Coreference Resolver

(Feb 2020)

The Classic Example:

- The monkey ate the banana because it was hungry.
 - What is it? Monkey or the banana?

MentionCoref--MentionMentionThe monkey ate the banana because it was hungry.

Slide credit: Omid Bakhshandeh

The paradigm shift in NLP, since 2015...

2015-2017:

- What happened: New SOTA established on various NLP benchmark
- *Recipe*: Encode the input text using BiLSTMs, decode with attention!
- Shortcomings: Could not tackle reading comprehension tasks that (supposedly) required:
 - Vast amount of background knowledge, or
 - Reasoning, or
 - Had long established contexts.
 - e.g., Story Cloze Test (Mostafazadeh et al., 2016).

3. The BiLSTM Hegemony

28

To a first approximation, the de facto consensus in NLP in 2017 is that no matter what the task, you throw a BiLSTM at it, with attention if you need information flow

Chris Manning

Story Cloze Test (Mostafazadeh et al., 2016) Narrative comprehension benchmark

Context:

Jim got his first credit card in college. He didn't have a job so he bought everything on his card. After he graduated he amounted a \$10,000 debt. Jim realized that he was foolish to spend so much money.

Two alternative endings:

Jim decided to devised a plan for repayment.

Jim decided to open another credit card.

A challenging commonsense reasoning task, where SOTA was ~65% for many months after release of the dataset.

Things got interesting in 2018!

- Late 2017-2018:
 - What happened. The dawn of Attention is All you need (Vaswani et al., 2017), introducing transformers. Brand new established SOTA on various supposedly more complex reading comprehension tasks.
 - *Recipe*: fine-tune large pretrained transformer-based models on downstream tasks (even with a small supervised data)!



Improving Language Understanding with Unsupervised Learning

We've obtained state-of-the-art results on a suite or diverse language tasks with a scalable, taskagnostic system, which we're also releasing. Our approach is a combination of two existing ideas: <u>transformers</u> and <u>unsupervised pre-training</u>. These results provide a convincing example that pairing supervised learning methods with unsupervised pre-training works very well; this is an idea that many have explored in the past, and we hope our result motivates further research into applying this

GPT-1 Model	
Radford et al,	201

DATASET

rd et al , 2018	result motivates further research into a idea on larger and more diverse datase	
TASK	8074	01105

ROCStories	Commonsense Reasoning	77.6	86.5
СОРА	Commonsense Reasoning	71.2	78.6

These results were on the Story Cloze Test v1, where there had been some stylistic biases (Sap et al., 2017). We tested a host of models on the new blind Story Cloze Test v 1.5 test set (Sharma et al., 2018). The GPT-1 model was the only model still holding its rather high performance!

So, are these models actually learning to transfer various lexical, conceptual, and world knowledge?

2019 was an exciting year for NEP!

- The 2018's recipe of transfer learning was impressively in full bloom in 2019! MI-DNN KD
- The community has started to think about the problems and weaknesses of the emerging techniques.



R. Thomas McCoy,¹ Ellie Pavlick,² & Tal Linzen¹ ¹Department of Cognitive Science, Johns Hopkins University ²Department of Computer Science, Brown University tom.mccoy@jhu.edu, ellie_pavlick@brown.edu, tal.linzen@jhu.edu

t is now almost a cliché to find out that BERT (Devlin et al., 2019) performs "surprisingly well" on whatever dataset you throw at it.

NLP's Clever Hans Moment has Arrived

Benjamin Heinzerling S

XLM

So have we come far enough?





Our moonshot at **Cognition**

Machines as thought partners!

We are working building Al systems that **build a shared understanding** with human and **explain** their answers well enough to eventually teach humans!



14

When humans, even young children, read, they make countless implicit commonsense inferences that frame their understanding of the unfolding narrative!



Peppa was riding her bike. A car turned in front of her. Peppa turned her bike sharply. She fell off of her bike. Peppa skinned her knee. While reading, humans construct a coherent representation of what happened and *why*, combining information from the text with relevant background knowledge

Humans can construct the causal chain that describes how the sequence of events led to a particular outcome!



A car turned in front of Peppa causes \rightarrow

Peppa to turn her bike sharply causes→

Peppa fell off of her bike

 $causes \rightarrow$

Peppa skinned her knee

 $causes \rightarrow$

(likely) she asks for help!

Humans can also describe how characters' different states, such as emotions and location, changed throughout the story.

Peppa was on her bike throughout riding it.

Then after falling, Peppa was on the ground.

eppa went from feeling (likely) happy to feeling in pain after falling.

Though humans build such mental models of situations with ease (Zwaan et al., 1995), **Al systems** for tasks such as reading comprehension and dialogue **remain far from exhibiting similar commonsense reasoning capabilities**!

Why?

• Two major bottlenecks in the AI research:

Not having ways of acquiring (often-implicit) commonsense knowledge at scale.

Not having ways to incorporate knowledge into the state-of-the-art AI systems.

GLUCOSE: GeneraLized and COntextualized Story Explanations!

A new commonsense reasoning framework for tackling both those bottlenecks at **scale**!





Jennifer Chu-Carroll

Lori Moon

Aditya Kalyanpur



Lauren Berkowitz

David Buchanan

GLUCOSE Commonsense Reasoning Framework

Given a short story S and a selected sentence X in the story, GLUCOSE defines ten dimensions of commonsense causal explanations related to X, inspired by human cognitive psychology.

GLUCOSE framework through an Example

Peppa was riding her bike. A car turned in front of her. Peppa turned her bike sharply She fell off of her bike. Peppa skinned her knee.

Semi-structured Inference Rule = antecedent *connective* consequent



GLUCOSE framework through an Example

Peppa was riding her bike. A car turned in front of her. Peppa turned her bike sharply. She fell off of her bike. Peppa skinned her knee.



Dim	Peppa po	ossesses	<u>a bike</u> <i>l</i>	Enables Peppa	turned l	her bike		
#4	subject	verb	object	subject	verb	object		
Is there a possession	Someon	e_A poss	esses Soi	$\underbrace{\text{nething}_A}_{\text{Ena}}$	bles Som	\underline{eone}_A <u>n</u>	noves	Something _A
state that enables X?	subject	ve	erb	object	su	bject	verb	object



N/A (the dimension is not applicable for this example)

Are there any other attributes enabling X?

GLUCOSE is a unique perspective on commonsense reasoning for presenting often-implicit commonsense knowledge in the form of *semi-structured general* inference rules that are also grounded in the context of a specific story!

GLUCOSE captures mini causal theories about the world focused around events, states (location, possession, emotion, etc), motivations, and naive human psychology.



How to address the problem of implicit knowledge acquisition at scale?

Filling in the GLUCOSE dimensions is **cognitively a complex task** for **lay workers**, since it **requires grasping the concepts of causality and generalization** and to write **semi-structured inference rules**!

An effective multistage crowdsourcing platform

After many rounds of pilot studies, we successfully designed an effective platform for collecting GLUCOSE data that is cognitively accessible to laypeople!

Read Instructions Frequently Asked Questions (FAQ Please answer the following queries about the story below **Quick MTurk Review Dashboard** The most thorough and accurate submissions will receive bonuses! We have many more HITs coming Account Balance: 6407 70 Story: Worker ID 3GR2F62BN49LR89U21N3JKD87TN Jennifer has a big exam tomorrow. She wants to nail the exam. She pulls an all-nighter. The next day, she is very tired. Her teacher tells the students that the test is postponed. Jennifer is quite relieved. 6F71BC7TJQE48V9ZTNN Let's call the highlighted sentence X = The next day, she is very tired An event that directly causes or enables X Now you can perform the following actions as you wish Consider the events that happen before X (or are likely to happen). Does any of them directly cause X, or simply make X possible (i.e., enable X)? Whenever possible, you are encouraged to find the answer from the other sentences in the story. Remember, there are often no right or wrong answers; just give us your intuition ATTAENY0EAN2SV4HKFNSZ 3GR2F62BN49LR89U21N3JKD87TN 36F71BC7TJQF48V9ZTNN1BQ84N Your Answer: No. I can't think of anything really/the query is not applicable to this sentence Yes, below is my two-step answer You can filter the list of submissions using the following fields: **GLUCOSE** Review Dashboard **GLUCOSE** Qualification UI You will see the submit button when you reach the end of the queries Thanks for your hard work! If you encounter any issues, please contact us. 18 GLUCOSE Main U

Statistics and Examples

Various implicit and script-like mini-theories:

- Someone_A gives Someone_B Something_A Results in Someone_B possess(es) Something_A
- Someone_A is Somewhere_A *Enables* Someone_A forgets Something_A Somewhere_A
- Someone_A is careless *Enables* Someone_A forgets Something_A Somewhere_A
- Someone_A forgets Something_A Somewhere_A Results in Something_A is Somewhere_A
- Someone_A feel(s) tired *Enables* Someone_A sleeps
- Someone_A is in bed *Enables* Someone_A sleeps
- Someone_A runs into Someone_B (who Someone_A has not seen for a long time) *Causes* Someone_A feel(s) surprised
- Someone_A asks Someone_B a question *Causes/Enables* Someone_B answers the question

# total inference rules	620K
# total unique stories	4700
# workers participated.	372
# mins per HIT on avg.	4.6mir

To our knowledge, GLUCOSE is among the few cognitivelychallenging AI tasks to have been successfully crowdsourced!



GLUCOSE captures extensive commonsense knowledge that is unavailable in the existing resources

Ceiling overlap between GLUCOSE and other resources based on besteffort mapping of relations.

GLUCOSE Dim	า1	2	5	6	7	10
ConceptNet	1.2%	0.3%	0%	1.9%	0%	0%
ATOMIC	7.8%	1.2%	2.9%	5.3%	1.8%	4.9%



How to incorporate commonsense knowledge into the state-of-the-art Al systems?

GLUCOSE Commonsense Reasoning Benchmark A testbed for evaluating models that can incorporate such commonsense knowledge and show inferential capabilities

- Task: Given a story *S*, the sentence <u>X</u>, and dimension *d*, the GLUCOSE specific and general rules should be predicted.
- Test Set: We carefully curated a doubly vetted test set, based on previously unseen stories and on which our most reliable annotators had high agreement. Our vetting process resulted in a test set of 500 GLUCOSE story/sentence pairs, each with 1-5 dimensions answered.
- Evaluation Metrics: Human and Automatic

Read Instruction

Frequently Asked Questions (FAC

Please rate the accuracy of various answers to the query about the story below. The most accurate ratings will receive bonuses!



Story:

Cody caught a mouse in his trap, He checked the trap after two weeks. He found the dead mouse. Cody threw the dead mouse in the trash. His cat dug the mouse out of the trash can.

Let's call the highlighted sentence X = He checked the trap after two weeks

An event that directly causes or enables X

Query: Consider the events that happen before X (or are likely to happen). Does any of them directly cause X, or simply make X possible (i.e., enable X)?

For each of the following candidate answers, taking into account the story, the highlighted sentence X, and your own common sense, rate the Specific Statement and/or General Rule on the scale of "incorrect" to "correct". Please take into account your own prior understanding of what a good Specific Statement or General Rule should look like in "Explain a Story" task. Following is what each rating means:

- Completely incorrect: This answer is completely irrelevant! Meaning, either it (a) is a completely irrelevant answer the above particular query about the particular selected sentence in the context of this particular story, and/or (2) has some major errors in how the content is composed that makes the answer incorrect.
- Almost Incorrect: This answer is "not" completely irrelevant, has some correct components with a few serious errors! Either it (a) is not really a correct answer for the
 above particular query, specially given the particular selected sentence in the context of this particular story, and/or (2) has a few notable error(s) in how the content is
- Almost Correct: This answer is overall correct but has some minor flaws. However, either it (1) is not a very accurate answer for the above particular query, specially given
 the particular selected sentence in the context of this particular story, and/or (2) has some minor error(s) in how the content is composed that make the answer a bit
 incoherent.
- Completely Correct: This answer is completely correct! Meaning, it is an accurate answer for the above particular query given the particular selected sentence in the
 context of this particular story.

Following are the candidate answers to the above question:

----- Candidate 1-----

Specific Statement Answer: the mouse was still alive >Causes/Enables> He checked the trap after two weeks.

Completely Incorrect O Almost Incorrect O Almost Correct O Completely Correct

----- Candidate 2 -----

Specific Statement Answer: Cody caught a mouse in his trap >Causes/Enables > Cody checked the trap after two weeks

Completely Incorrect Almost Incorrect Almost Correct Completely Correct

----- Candidate 3 -----

Specific Statement Answer: Cody catches a mouse in his trap >Causes/Enables > Cody checks the trap after two weeks

Completely Incorrect Almost Incorrect Completely Correc
 Completely Correc

Specific Statement Answer: Cody caught a mouse >Causes/Enables> He checked the trap after two weeks.

Completely Incorrect Completely Correct Completely Correct

----- Candidate 5 -----

Specific Statement Answer: a cat eats his shoe >Causes/Enables> he checks the trap

Completely Incorrect
Almost Incorrect
Completely Correct
Completely C

Specific Statement Answer: Cody sets a trap to catch a mouse >Causes/Enables> Cody checks the trap

Completely Incorrect O Almost Incorrect O Almost Correct O Completely Correct

We designed a specialized Human Evaluation UI for collecting reliable, reproducible, and calibrated ratings!

Automatic Evaluation

of natural language generations

- A majority of commonsense reasoning frameworks have been in multiple-choice form, as opposed to natural language generation, due to ease of evaluation
 - Multiple-choice tests are inherently easier to be **gamed**!
- Automatic evaluation for tasks involving natural language generation with diverse possibilities has been a major bottleneck for research
- BLEU's ease of replicability has made it a popular automated metric, but its correlation with human judgement has proven weak in various tasks.

Automatic Evaluation

of natural language generations in GLUCOSE

- We found very strong pairwise correlation between human and ScareBLEU corpus-level scores on our test set.
 - Spearman = 0.891, Pearson = 0.855, and Kendall's = 0.705, all with p-value < 0.001.
- This is accomplished through various design choices in GLUCOSE:
 - 1) GLUCOSE **semi-structured inference rules** are **designed to be evaluable**, where the **structure naturally limits the format** of the generated rules
 - 2) We curated our test set to eliminate cases with a wide range of correct responses where humans cannot agree, making the limited number of gold references sufficient for automatic evaluation
 - 3) We designed a systematic human evaluation process that can **collect calibrated ratings from judges** who are well educated about what constitutes a correct GLUCOSE rule.



GLUCOSE task has a systematic evaluation that is fast and easily replicable!

Notable Models & Results

Baselines KNN GPT-2 Avg: g 0.5/3 Avg: s 0.8/3

We show that:

1) The KNN model performs the worst, highlighting the importance of generalizing beyond the training data.

2) Pre-trained language model perform very poorly at the task and does not show basic commonsense inference!

3) When the pre-trained neural models are fine-tuned on the rich GLUCOSE data, they achieve very high performance in making commonsense predictions on *unseen* stories.



Example Predictions

Dimension 3; a location enabling X.

- Input:
 - Karen made a pan of lasagna. She brought it to the party. Nobody wanted to eat lasagna. Karen ate it for a week. She became tired of lasagna.

GPT- 2	She w N/A	as in front of a TV >Enables> Karen made a pan of lasagna.
Full-LM		Karen is at home >Enables> Karen made a pan of lasagna and brought it to the party Someone _A is in Somewhere _A >Enables> Someone _A makes Something _A (that is edible)
Enc-Dec		Karen is in the kitchen > <i>Enables</i> > Karen makes a pan of lasagna Someone _A is in a kitchen >Enables> Someone _A cooks Something _A
Human		Karen is in the kitchen > <i>Enables</i> > Karen made a pan of lasagna Someone _A is in a kitchen >Enables> Someone _A prepares Something _A (that is a dish)

Example Predictions

Dimension 6; an event that <u>X</u> Causes/Enables.

- Input:
 - Karen made a pan of lasagna. She brought it to the party. Nobody wanted to eat lasagna. Karen ate it for a week. She became tired of lasagna.



We proved our following hypothesis

A promising new recipe for giving machines commonsense is to use **high-quality commonsense knowledge** as the **seed data** for training **neural models that have pre-existing lexical and conceptual knowledge**.

Static commonsense knowledge base with GLUCOSE mini-theories authored by humans



GLUCOSE-Trained model that can generate GLUCOSE dimensions for any novel input

Old-school commonsense knowledge base is static

Modern commonsense knowledge base is dynamic

To conclude:

We've come a rather long way in the last decade in NLP with lots of exciting progress. My hope for our directions in 2020 is to work on tackling the following issues which we are still grappling with...

- Our amazing models sometimes make glaringly stupid mistakes, being brittle! This makes it hard to deploy these models into real-world products.
- We yet don't know the implications of establishing SOTA on various benchmarks. Are we
 making any real progress? Do these models work outside of our narrow lab settings in the
 real world?
- We still cannot tackle tasks that have little to no annotated data. Better knowledge transfer across domains and incorporating prior knowledge and world models is essential.
- Handful of top industry players get to pay the costs for building ever-larger (and not-green) models. Where are we going with this paradigm?
- And we yet don't have an AI system that has commonsense of perhaps even a dog (?), let alone a 5-year-old kid...

Thanks for listening!



31