

Infusing Disease Knowledge into BERT for Health Question Answering, Medical Inference and Disease Name Recognition

Yun He¹, Ziwei Zhu¹, Yin Zhang¹, Qin Chen², James Caverlee¹ ¹Texas A&M University, College Station, USA ²Fudan University, Shanghai, China



- Disease knowledge is critical for many health and biomedical related tasks such as consumer health question answering.
- In this paper, we integrate BERT models with disease knowledge for improving these important tasks.
- Our proposed method, Disease Knowledge Infusion Training, achieves new SOTA results in two datasets: MEDIQA-2019 and MEDNLI.

Disease

• Disease has a huge impact on human life.



TEXAS A&M

ĀM

• Disease is one of the fundamental biological entities in biomedical research.

Until October 12, 2020, COVID-19 has killed 1.08M people and disrupted the everyday lives of billions of people worldwide.

Disease Knowledge

Knowledge of a disease Fever is the most common symptom, includes information of but highly variable in severity and Symptoms presentation, with some older... various aspects of the The standard method of testing is realdisease. Diagnosis time reverse transcription polymerase chain reaction (rRT-PCR)... • This disease knowledge Knowledge of Covid-19 is critical for many People are managed with supportive care, which may include fluid therapy, Treatment health and biomedical oxygen support, and supporting... related NLP tasks.

TEXAS A&M

ĀM

Consumer Health Question Answering

Question: ...keen to learn how to get COVID-19 diagnosed, many thanks

Answer 1: ... real-time reverse transcription polymerase chain reaction... Answer 2: ... diagnosis of vipoma requires demonstration of diarrhea... Answer 3: ...affected by this disorder are not able to make lipoproteins...

Label: Answer 1 is the most relevant Disease Knowledge: Answer 1 is the diagnosis of COVID-19



TEXAS A&M

Medical Language Inference

Premise: She was not able to speak, but appeared to comprehend well

Hypothesis: Patient had aphasia

Label: entailment Disease Knowledge: Premise describes the symptoms of aphasia



I FXA

ĀŇ

Disease Name Recognition

Text: **Myotonic dystrophy** (DM) is **caused by a CTG expansion** in the 3 untranslated region of the DM gene.

IEXA

Ā M

Label: Myotonic dystrophy

Disease Knowledge: the text contains the cause of Myotonic dystrophy



Infusing Disease Knowledge into BERT models





I want to know the semantic relations among them so that I can tell which disease and aspect is this diseasedescriptive text talking about.

ĀM

TEXAS A&M

Infusing Disease Knowledge into BERT models



TEXAS A&M

Ă Ň

Resources of Disease Terms

First, we seek a disease vocabulary that provides disease terms. Several resources include:







SNOMED CT The global language of healthcare

TEXAS A&M

As a first step, we choose MeSH, which is a comprehensive controlled vocabulary proposed by the National Library of Medicine. Infections;C01 Aneurysm, Infected;C01.069 Arthritis, Infectious;C01.100 Arthritis, Reactive;C01.100.500 Asymptomatic Infections;C01.125 Bacterial Infections and Mycoses;C01.150 Bacterial Infections;C01.150.252 Bacteremia;C01.150.252.100

ĀΜ

Examples of MeSH Disease Terms

Proposed Method: Disease Knowledge Infusion Training

1. Obtain disease terms from MeSH



WikipediA The Free Encyclopedia

2. Obtain Articles of diseases from Wikipedia⁶ Management

COVID-19

Signs and symptoms

- 1.1 Complications
- 2 Cause
 - 2.1 Transmission
 - 2.2 Virology
- 3 Pathophysiology
 - 3.1 Immunopathology

4 Diagnosis

- 4.1 Pathology
- 5 Prevention

Passage: The standard method of testing is real-time reverse transcription polymerase chain reaction (rRT-PCR)...

Disease: COVID-19 (title of the Wikipedia article)

Ā M

Aspect: Diagnosis (title of the section)

4. Extract the weaklysupervised topic

disease and aspect.

TEXAS A&M

3. Extract text from a

section as the passage.

Auxiliary Sentence: What is the **diagnosis** of **COVID-19**? **5.** Construct an auxiliary sentence that mentions the topic disease and aspect.

New Passage for MLM:

What is the **[MASK]** of **[MASK]**? The standard method of testing is real-time reverse transcription polymerase chain reaction (rRT-PCR)...

6. Concatenate the passage and the auxiliary Sentence. **BERT** is trained to infer the disease and aspect.

Experiments

We infuse disease knowledge into these BERT models:

TEXAS A

Ā Ň

- General BERT
 - BERT-base (Wikipedia, BookCorpus)
 - ALBERT (Wikipedia, BookCorpus)
- Biomedical BERT
 - BioBERT (PubMed articles)
 - ClinicalBERT (Clinical notes)
 - BlueBERT (PubMed abstracts and clinical notes)
 - SciBERT (18% of papers from the computer science domain and 82% from the biomedical domain)

Relationship between Our Approach and other BERT Models



Continual Pre-training BERT with MLM and NSP on biomedical corpora.



Continual Pre-training BioBERT with proposed Disease Knowledge Infusion Training.



BioBERT + Disease Knowledge



General BERT:

Pre-training on Wikipedia

from scratch

WIKIPEDIA The Free Encyclopedia









WIKIPEDIA The Free Encyclopedia

Relationship with Biomedical BERT

 The biomedical BERT models (e.g., ClinicalBERT and BioBERT) capture the general syntactic and semantic knowledge of biomedical language.

ĀM

- Our method is specifically designed for capturing the semantic relations between a disease-descriptive text and its corresponding aspect and disease.
- These biomedical BERT models can be further enhanced by our approach in biomedical NLP tasks.

Consumer Health Question Answering

Our approach can enhance BERT models in nearly all cases. New SOTA results are obtained in MEDIQA-2019 datasets.

Accuracy (%) on MEDIQA-2019 Dataset

TEXAS A&M



Disease Name Recognition

Our approach does not outperform SOTA but it can enhance BERT models in most cases.

TEXAS A&M

ĀŇ



Our approach can enhance BERT models in nearly all cases. New SOTA results are obtained in MEDNLI datasets.

TEXAS A&M

Ă Ň



Conclusion

- We propose Disease Knowledge Infusion Training to augment BERT-like models with disease knowledge.
- Our approach can enhance both general BERT models (e.g., BERT and ALBERT) and biomedical BERT models (e.g., ClinicalBERT and BioBERT).

IEXA

 New SOTA results are obtained in MEDIQA-2019 and MEDNLI datasets.



Implicative Reasoning in Conversational Machine Reading

Yifan Gao, The Chinese University of Hong Kong Nov 4, 2020

Machine Comprehension

SQuAD: 100,000+ Questions for Machine Comprehension of Text (Rajpurkar et al., 2016)

The current Chief Executive is **Carrie Lam**, who was selected on 26 March 2017, appointed by the Central People's Government with the State Council Decree signed by Premier Li Keqiang, on 11 April 2017 and took office on 1 July 2017.

✓ Literal Answer



Q: Who is the chief executive of Hong Kong?



Conversational Question Answering

CoQA: A Conversational Question Answering Challenge (Reddy et al., 2018)

Incumbent **Democratic** President Bill Clinton was ineligible to serve a **third term** due to **term limitations** in the 22nd Amendment of the Constitution, and Vice President Gore was able to secure the Democratic nomination with relative ease.

- ✓ Literal Answer
- Dialog Understanding



Q: What political party is Clinton a member of?

Q: Wha

Q: What was he ineligible to serve?







A: Democratic

However...

Interpreting Natural Language Rules

The text to read may <u>not</u> contains the literal answer, but it contains a **recipe** to derive it.



Scenario: I am a 34-year-old man from the United States who owns their own business. We are an American small business.

00

Question: Is the 7(a) Loan Program for me?

Implicative Reasoning in Conversational Machine Reading Introduction



Scenario: I am a 34-year-old man from the United States who owns their own business. We are an American small business.



Question: Is the 7(a) Loan Program for me?



Rule Text

7(a) loans are the most basic and most used type loan of the Small Business Administration's (SBA) business loan programs. Its name comes from section 7(a) of the Small Business Act, which authorizes the agency to provide business loans to American small businesses. The loan program is designed to assist for-profit businesses that are not able to get other financing from other resources.

Implicative Reasoning in Conversational Machine Reading Introduction



Task Definition

ShARC: Shaping Answers with Rules through Conversation



(Saeidi et al., 2018)



Rule Text: 7(a) loans are the most basic and ...

Scenario: I am a 34-year-old man

Question: Is the 7(a) Loan Program for me?

Dialog History

Follow-up Q1: Are you a for-profit business? A1: Yes

Decision Making

Make a prediction among:

Yes, No, Irrelevant, Inquire

- Yes/No: Directly answer the question
- Irrelevant: unanswerable
- Inquire-

Question Generation

Ask a follow-up question to clarify the unknown user information

Outline

1. Explicit Memory Tracker with Coarse-to-Fine Reasoning for Conversational Machine Reading, ACL 2020

2. Discern: Discourse-Aware Entailment Reasoning Network for Conversational Machine Reading, EMNLP 2020

The 58th Annual Meeting of the Association for Computational Linguistics

July 5 – 10 2020

Explicit Memory Tracker with Coarse-to-Fine Reasoning for Conversational Machine Reading

Yifan Gao¹, Chien-Sheng Wu², Shafiq Joty², Caiming Xiong², Richard Socher², Irwin King¹, Michael R. Lyu¹, Steven C.H. Hoi²

1. The Chinese University of Hong Kong 2. Salesforce Research

Code & Models: https://github.com/Yifan-Gao/explicit_memory_tracker

Contributions

- Explicit Memory Tracker (EMT)
 - Explicitly track whether conditions listed in the rule text have been fulfilled or not
- * Coarse-to-fine (C2F) Reasoning
 - A coarse-to-fine approach to reason out which part of the rule text is underspecified, and ask a question accordingly
- Our proposed solution achieves new state-of-the-art results on the ShARC benchmark





ACL 2020, Explicit Memory Tracker with Coarse-to-Fine Reasoning for Conversational Machine Reading

Encoding



- 1. Parse the rule text into multiple rule sentences according to rules
- 2. Insert **[CLS]** token at the start of each rule sentence, initial question, scenario, and question-answer pairs in the dialog history
- 3. Concatenate all information and feed to BERT for encoding
- 4. **[CLS]** symbol is treated as the feature representation of the sentence that follows it

Explicit Memory Tracking

Rule sentences $\mathbf{k}_1, \mathbf{k}_2, \dots, \mathbf{k}_M$



User provided information:

- Initial question \mathbf{s}_Q
- **\diamond** Scenario s_S
- Dialog history $\mathbf{s}_1, \dots, \mathbf{s}_P$



- We propose Explicit Memory Tracker (EMT), a gated recurrent memory-augmented neural network
- EMT explicitly tracks the states of rule sentences by sequentially reading the user provided information

Explicit Memory Tracking

EMT assigns a state \mathbf{v}_i to each key \mathbf{k}_i , and sequentially reads user information

At time step *t*:

$$\begin{split} \tilde{\mathbf{v}}_{i,t} &= \operatorname{ReLU}(\mathbf{W}_k \mathbf{k}_i + \mathbf{W}_v \mathbf{v}_{i,t} + \mathbf{W}_s \mathbf{s}_t), \\ g_i &= \sigma(\mathbf{s}_t^\top \mathbf{k}_i + \mathbf{s}_t^\top \mathbf{v}_{i,t}) \in [0, 1], \\ \mathbf{v}_{i,t} &= \mathbf{v}_{i,t} + g_i \odot \tilde{\mathbf{v}}_{i,t} \in \mathbb{R}^d, \mathbf{v}_{i,t} = \frac{\mathbf{v}_{i,t}}{\|\mathbf{v}_{i,t}\|} \end{split}$$



Keys and final states of rule sentences are denoted as $(\mathbf{k}_1, \mathbf{v}_1), \dots, (\mathbf{k}_M, \mathbf{v}_M)$

- Decision Making Module
- Question Generation Module

Proposed Solution Decision Making

Based on the most up-to-date key-value states of rule sentences $(\mathbf{k}_1, \mathbf{v}_1), \dots, (\mathbf{k}_M, \mathbf{v}_M)$, EMT makes a decision among *Yes, No, Irrelevant, Inquire*

$$egin{aligned} &lpha_i = \mathbf{w}_{lpha}^{ op}[\mathbf{k}_i;\mathbf{v}_i] + b_{lpha} \in \mathbb{R}^1 \ & ilde{lpha}_i = ext{softmax}(lpha)_i \in [0,1] \ & extbf{c} = \sum_i ilde{lpha}_i[\mathbf{k}_i;\mathbf{v}_i] \in \mathbb{R}^d \ & extbf{z} = \mathbf{W}_z \mathbf{c} + \mathbf{b}_z \in \mathbb{R}^4 \end{aligned}$$



The decision making module is trained with the following cross entropy loss:

$$\mathcal{L}_{dec} = -\log \operatorname{softmax}(\mathbf{z})_l$$

Subtask: Entailment State Prediction

- Explicitly track whether a condition listed in the rule has already been satisfied or not
- ✤ The possible entailment labels are:
 - Entailment (E)
 - Contradiction (C)
 - Unknown (U)

$$\mathbf{e}_{i} = \mathbf{W}_{e}[\mathbf{k}_{i}; \mathbf{v}_{i}] + \mathbf{b}_{e} \in \mathbb{R}^{3}$$
$$\mathcal{L}_{\text{entail}} = -\frac{1}{M} \sum_{i=1}^{M} \log \operatorname{softmax}(\mathbf{e}_{i})_{r}$$



Follow-up Question Generation

When the decision is 'Inquire', a <u>follow-up question</u> is required for further clarification.

We adopt a two-step approach:

- 1. Extract a span inside the rule text which contains the underspecified user information
- 2. Rephrase the extracted span into a follow-up question
Follow-up Question Generation: Coarse-to-fine Underspecified Span Extraction

- 1. Coarse-to-fine Underspecified Span Extraction
 - 1) Identify underspecified rule sentence ζ_i



ACL 2020, Explicit Memory Tracker with Coarse-to-Fine Reasoning for Conversational Machine Reading

Follow-up Question Generation: Coarse-to-fine Underspecified Span Extraction

- 1. Coarse-to-fine Underspecified Span Extraction
 - 1) Identify underspecified rule sentence ζ_i
 - 2) Extract a span within each rule sentence $(\gamma_{i,j}, \delta_{i,j})$



ACL 2020, Explicit Memory Tracker with Coarse-to-Fine Reasoning for Conversational Machine Reading

Follow-up Question Generation: Coarse-to-fine Underspecified Span Extraction

- 1. Coarse-to-fine Underspecified Span Extraction
 - 1) Identify underspecified rule sentence ζ_i
 - 2) Extract a span within each rule sentence $(\gamma_{i,j}, \delta_{i,j})$
 - 3) Select the span with the highest span score $\zeta_i * (\gamma_{i,j}, \delta_{i,j})$



Follow-up Question Generation

2. Question Rephrasing

- 1) Finetune UniLM (Dong et al, 2019), a pretrained language model
- 2) [CLS] rule text [SEP] span [SEP]



Overall Loss for EMT

The overall loss is the sum of the decision loss, entailment prediction loss and span extraction loss:

$$\mathcal{L} = \mathcal{L}_{dec} + \lambda_1 \mathcal{L}_{entail} + \lambda_2 \mathcal{L}_{span}$$

Experimental Setup

✤ Dataset:

- ShARC CMR dataset (Saeidi et ¹
- Train/Dev/Test dataset sizes are
- Test set is not public.
- Leaderboard: <u>https://sharc-data.</u>

ShARC: End-to-end Task

#	Model / Reference	Affiliation	Date	Micro Accuracy[%]	Macro Accuracy[%]	BLEU-1	BLEU-4
1	[Anonymous]	[Anonymous]	May 2020	73.2	78.3	64.0	49.1
2	EMT	Salesforce Research & CUHK	Nov 2019	69.4	74.8	60.9	46.0
3	EMT + entailment	Salesforce Research & CUHK	Mar 2020	69.1	74.6	63.9	49.5
4	[Anonymous]	[Anonymous]	Dec 2019	69.0	74.6	56.7	42.0
5	E3	University of Washington	Feb 2019	67.6	73.3	54.1	38.7
6	BiSon (single model)	NEC Laboratories Europe	Aug 2019	66.9	71.6	58.8	44.3

Experimental Setup

ShARC: End-to-end Task

						Micro	Macro		
•	Dataset	#	Model / Reference	Affiliation	Date	Accuracy[%]	Accuracy[%]	BLEU-1	BLEU-4
•	 ShARC CMR dataset (Saeidi et Train/Dev/Test dataset sizes are 	1	[Anonymous]	[Anonymous]	May 2020	73.2	78.3	64.0	49.1
	 Test set is not public. Leaderboard: <u>https://sharc-data.</u> 	2	EMT	Salesforce Research & CUHK	Nov 2019	69.4	74.8	60.9	46.0
**	Evaluation MetricsEnd-to-End Evaluation	3	EMT + entailment	Salesforce Research & CUHK	Mar 2020	69.1	74.6	63.9	49.5
		٨		[Apopymous]	Dec	69.0	74.6	567	42.0

If two models have different *Inquire* predictions, the follow-up questions for evaluation will be different, making the comparison unfair.

6	BiSon (single	NEC Laboratories	Aug	66.9	71.6	58.8	44.3
	model)	Europe	2019				

Experimental Setup

- ✤ Dataset:
 - ShARC CMR dataset (Saeidi et al. 2018)
 - Train/Dev/Test dataset sizes are 21980/2270/8276.
 - Test set is not public.
 - Leaderboard: <u>https://sharc-data.github.io/leaderboard.html</u>
- Evaluation Metrics
 - End-to-End Evaluation
 - Oracle Question Generation Evaluation
 - We propose a new evaluation perspective.
 - We ask the models to generate follow-up questions *whenever* the ground truth decision is <u>Inquire</u>, and compute the BLEU score.

Leaderboard Submission

Models	End-to-End Task (Leaderboard Performance)				
WIOUEIS	Micro Acc.	Macro Acc.	BLEU1	BLEU4	
Seq2Seq (Saeidi et al., 2018)	44.8	42.8	34.0	7.8	
Pipeline (Saeidi et al., 2018)	61.9	68.9	54.4	34.4	
BERTQA (Zhong and Zettlemoyer, 2019)	63.6	70.8	46.2	36.3	
UrcaNet (Sharma et al., 2019)	65.1	71.2	60.5	46.1	
BiSon (Lawrence et al., 2019)	66.9	71.6	58.8	44.3	
E ³ (Zhong and Zettlemoyer, 2019)	67.6	73.3	54.1	38.7	
EMT (our single model)	69.1	74.6	63.9	49.5	

Table 1: Performance on the blind, held-out test set of ShARC end-to-end task.

Class-wise Decision Prediction Accuracy

Models	Yes	No	Inquire	Irrelevant
BERTQA	61.2	61.0	62.6	96.4
E^3	65.9	70.6	60.5	96.4
UrcaNet*	63.3	68.4	58.9	95.7
EMT	70.5	73.2	70.8	98. 6

Table 2: Class-wise decision prediction accuracy on the development set (*: reported in the paper).

Oracle Question Generation Task

	Oracle Question Generation Task					
Models	Develop	ment Set	Cross Va	alidation		
	BLEU1	BLEU4	BLEU1	BLEU4		
$-E^3$	52.79±2.87	37.31±2.35	51.75	35.94		
E ³ +UniLM	57.09±1.70	$41.05 {\pm} 1.80$	56.94	42.87		
EMT	62.32 ±1.62	47.89 ±1.58	64.48	52.40		

Table 3: Performance on Oracle Question Generation Task. We show both results on the development set and 10-fold cross validation. E^3 +UniLM replaces the editor of E^3 to our finetuned UniLM.

Experiment Interpretability	E : Entailment C : Contradicti	on U	: Unkno	wn	β_{unknown}		
	Regulation Text A	Entailment States					
	(parsed into six rule sentences: S1 \sim S6)			Turn 3			
	S1 Statutory Maternity Pay			U (99.99)			
	S2 To qualify for smp you must:	U (99.99)	U (99.99)	U (99.99)			
	S3 * earn on average at least £113 a week	U (99.93)	E (99.91)	E (99.67)			
	S4 * give the correct notice	U (99.97)	U (99.61)	C (99.81)			
	S5 * give proof you're pregnant	U (99.98)	U (99.75)	U (99.94)			
	S6 * have worked for your employer	U (99.98)	U (99.70)	U (99.96)			
	Scenario: I've been old enough to get my pension. My wife just reached pension age last year. Neither of us have applied for it yet.						
	Decision: Generated Question			Answer			
Turn 1: Inquire Do you earn on average at least £113 a wee				Yes			
	Turn 2: Inquire Did you give the correct notice?						
Γ	Turn 3: No						

ACL 2020, Explicit Memory Tracker with Coarse-to-Fine Reasoning for Conversational Machine Reading

Conclusion

- We propose a new approach called Explicit Memory Tracker (EMT) for conversational machine reading.
- EMT achieved a new state-of-the-art result on the ShARC CMR challenge.
- * EMT also gains interpretability by showing the entailmentoriented reasoning process as the conversation flows.

EMNLP 2020

The 2020 Conference on Empirical Methods in Natural Language Processing

Discern: Discourse-Aware Entailment Reasoning Network for Conversational Machine Reading

Yifan Gao¹ Chien-Sheng Wu² Jingjing Li¹ Shafiq Joty^{2,3}
Steven C.H. Hoi² Caiming Xiong² Irwin King¹ Michael R. Lyu¹
1. The Chinese University of Hong Kong
2. Salesforce Research 3. Nanyang Technological University

Code & Models: https://github.com/Yifan-Gao/Discern

- 1. Document Interpretation
 - Identification of Conditions
 - Determination of Logical Structures

7(a) loans are the most basic and most used type loan of the Small Business Administration's (SBA) business loan programs. Its name comes from section 7(a) of the Small Business Act, which authorizes the agency to provide business loans to American small businesses. The loan program is designed to assist for-profit businesses that are not able to get other financing from other resources.

"Eligible for 7(a) loans" = (1 = True) and (2 = True) and (3 = True)

- 2. Dialog Understanding
 - Track the user's fulfillment over the conditions
 - Jointly consider the fulfillment states and the logical structure of rules

7(a) loans are the ... provide business loans to American small businesses. The loan program is designed to assist for-profit businesses that are not able to get other financing from other resources.

- 2. Dialog Understanding
 - Track the user's fulfillment over the conditions
 - Jointly consider the fulfillment states and the logical structure of rules

7(a) loans are the ... provide business loans to American small businesses. The loan program is designed to assist for-profit businesses that are not able to get other financing from other resources.

- 2. Dialog Understanding
 - Track the user's fulfillment over the conditions
 - Jointly consider the fulfillment states and the logical structure of rules

7(a) loans are the ... provide business loans to American small businesses. The loan program is designed to assist for-profit businesses that are not able to get other financing from other resources.

Fulfillment State: (1)==True) and (2) == True) and (3) == True \Rightarrow Decision: Yes (You can apply the loan.)

Overview



Overview



Rule Segmentation

• Goal

- Understand the logical structure of the rule text
- Parse the rule into individual conditions for entailment reasoning
- Challenges
 - Sentence splitting is not enough: one rule sentence may contain several in-line conditions
- Solution: Discourse Segmentation
 - In the Rhetorical Structure Theory (RST) of discourse parsing (Mann and Thompson, 1988), texts are split into clause-like units called <u>elementary discourse units (EDUs)</u>

Rule Text: If a worker has taken more leave than they're entitled to, their employer must not take money from their final pay unless it's been agreed beforehand in writing.

Discourse Segmentation

\downarrow

[If a worker has taken more leave than they're entitled to,]_{EDU1} [their employer must not take money from their final pay]_{EDU2} [unless it's been agreed beforehand in writing.]_{EDU3}

Overview



Discern: Decision Making via Entailment Reasoning

Encoding



Discern: Decision Making via Entailment Reasoning

Entailment Prediction



Multi-Sentence Entailment Prediction

Discern: Decision Making via Entailment Reasoning

Entailment Prediction



Overview



Decision Making



EMNLP 2020, Discern: Discourse-Aware Entailment Reasoning Network for Conversational Machine Reading

Main Results

Models	End-to-End Task (Leaderboard Performance)				
	Micro Acc.	Macro Acc.	BLEU1	BLEU4	
Seq2Seq (Saeidi et al., 2018)	44.8	42.8	34.0	7.8	
Pipeline (Saeidi et al., 2018)	61.9	68.9	54.4	34.4	
BERTQA (Zhong and Zettlemoyer, 2019)	63.6	70.8	46.2	36.3	
UrcaNet (Sharma et al., 2019)	65.1	71.2	60.5	46.1	
BiSon (Lawrence et al., 2019)	66.9	71.6	58.8	44.3	
E^3 (Zhong and Zettlemoyer, 2019)	67.6	73.3	54.1	38.7	
EMT (Gao et al., 2020)	69.4	74.8	60.9	46.0	
EMT+entailment (Gao et al., 2020)	69.1	74.6	63.9	49.5	
DISCERN (our single model)	73.2	78.3	64.0	49.1	

Ablation Study

Models	Micro Acc.	Macro Acc.	
DISCERN	74.97 ± 0.27	$79.55 {\pm} 0.35$	_
DISCERN (BERT)	73.07 ± 0.21	$77.77 {\pm} 0.24$	RoBERTa > BERT
DISCERN (w/o EDU)	73.34 ± 0.22	$78.25{\scriptstyle\pm0.57}$	Discourse Segmentation > Sentence Splitting
DISCERN (w/o Trans)	74.25 ± 0.36	$78.78{\scriptstyle\pm0.57}$	Inter-Sentence Transformer IS Necessary!
DISCERN (w/o $\mathbf{\tilde{e}}$)	73.55 ± 0.26	$78.19{\pm}0.30$	Both Condition Representations and
DISCERN (w/o \mathbf{V}_{EDU})	72.95 ± 0.23	$77.53{\scriptstyle\pm0.19}$	Sentailment Vectors Facilitate Decisions

Analysis of Logical Structure of Rules



How Far Has the Problem Been Solved?

Idea: Disentangle the challenge between scenario interpretation and dialog understanding

- Full Dataset: Full development set of ShARC
- **Dialog Understanding Subset**: User information only contains dialog history
- Scenario Interpretation Subset: User information only contains user scenario



EMNLP 2020, Discern: Discourse-Aware Entailment Reasoning Network for Conversational Machine Reading

Conclusion

- We present Discern, a system that does discourse-aware entailment reasoning for conversational machine reading.
- Results on the ShARC benchmark shows that Discern outperforms existing methods by a large margin.
- We also conduct comprehensive analyses to unveil the limitations of Discern and challenges for ShARC.

Thanks!



Keep CALM and Explore: Language Models for Action Generation in Text-based Games

> <u>Shunyu Yao</u>, Rohan Rao, Matthew Hausknecht, Karthik Narasimhan

> > **EMNLP 2020**

RL view



- Discrete MDP
- Millions of states
- Billions of actions
- Very sparse reward

Observation: West of House You are standing in an open field west of a white house, with a boarded front door. There is a small mailbox here. Action: **Open mailbox** Observation: Opening the small mailbox reveals a leaflet. Action: Read leaflet Observation: (Taken) "WELCOME TO ZORK! ZORK is a game of adventure, danger, and low cunning. In it you will explore some of the most amazing territory ever seen by mortals. No computer should be without one!" Action: Go north Observation: North of House You are facing the north side of a white house. There is no door here, and all the windows are boarded up. To the north a narrow path winds through the trees.

Human view



- Text adventure
- Places and rooms
- Some valid actions
- Commonsense

Difference: Language understanding!

Text-based Game

- A partially observable Markov decision process (POMDP) where a player
 - Issues text actions
 - Receives text observations and scalar rewards

• Immediate challenge: action space

- Random actions: billions, mostly not admissible
- Previous models: heuristic rules or action handicap
 - Problem: ad-hoc, hard to generalize
- Our approach: train language models for action generation
 - Scalable, generalizable

Observation: You are in the living room. There is a doorway to the east, a wooden door with strange gothic lettering to the west, which appears to be nailed shut, a trophy case, and a large oriental rug in the center of the room. You are carrying: A brass lantern ...

Random Actions:

close door north a eat troll with egg, ...
CALM (n-gram) Actions:
enter room, leave room, lock room,
open door, close door, knock on door, ...
CALM (GPT-2) Actions:

east, open case, get rug, turn on lantern, move rug, unlock case with key, ...

Next Observation: With a great effort, the rug is moved to one side of the room, revealing the dusty cover of a closed trap door...
Pontiextsialpartiechlesinguage Model (CALM)



Action candidates: handicap \rightarrow language model

1. Fair (no forbidden functions used) 2. Generalize (to unsupported games)

DRRN: He, et al. "Deep Reinforcement Learning with a Natural Language Action Space."

How to train CALM

 o_t : observation at time t a_t : action at time t

- Train CALM on human gameplay trajectories of form
 (o₁,a₁,o₂,a₂,...), predict action a_t based on context c_t=(o_{t-1},a_{t-1},o_t)
- Two language models:
 - 1. GPT-2: standard cross entropy loss on action tokens

$$\mathcal{L}_{\text{LM}}(\theta) = -\mathbb{E}_{(a,c)\sim D} \log p_{\theta}(a|c) \qquad p_{\theta}(a|c) = \prod_{i=1}^{m} p_{\theta}(a^{i}|a^{$$

- 2. N-gram:
 - Probability: only base on action n-grams, independent of context
 - Generation: only consider actions with objects seen in observation

How to use CALM in RL

- Given context c, use **trained**, **frozen** CALM to produce top-k candidates (k=30)
- RL as action re-ranker to maximize game performance
- Rationale: combine CALM's generic action priors and RL's ability to optimize gameplay performance

Observation: You are in the living room. There is a doorway to the east, a wooden door with strange gothic lettering to the west, which appears to be nailed shut, a trophy case, and a large oriental rug in the center of the room. You are carrying: A brass lantern ...

Random Actions:

close door, north a, eat troll with egg, ... CALM (n-gram) Actions:

enter room, leave room, lock room, open door, close door, knock on door, ... CALM (GPT-2) Actions:

east, open case, get rug, turn on lantern, move rug, unlock case with key, ...

Next Observation: With a great effort, the rug is moved to one side of the room, revealing the dusty cover of a closed trap door...

		ClubFloyd Dataset	Jericho Walkthroughs
Data	# unique games	590	28
	Vocab size	39,670	9,623
	vocab size (game avg.) Ava trajectory length	2,303	1,057
	Action Quality	Non-optimal	Optimal

ClubFloyd Dataset

- Used for CALM training
- 590 games, 223,527 samples
- Human actions: noisy, nonoptimal, sometimes invalid
- A few very common actions (north, take all, examine, ...)

Jericho Benchmark

- Used to evaluate CALM + RL
- 28 games with simulator access
- Admissible action handicap at each step
- <u>Optimal walkthrough</u> for each game

No overlapping games!

Evaluate CALM with RL

- Baselines:
 - NAIL: ad-hoc rules, no RL or handicap
 - **DRRN**: need handicap for action space
 - KG-A2C: need handicap for action generation supervision, use KG
- No handicap: GPT-2 > N-gram ~ NAIL
- Handicap models: GPT-2 outperforms both DRRN and KG-A2C on 8/28 games
 - No "stupid" admissible actions
 - Handicap might miss some valid actions

NAIL: Hausknecht, et al. "Nail: A general interactive fiction agent."

KG-A2C: Ammanabrolu, Hausknecht. "Graph constrained reinforcement learning for natural language action spaces."

Avg. Norm. Score		
CALM (GPT-2)	9.4%	
CALM (n-gram)	5.5%	
NAIL	5.6%	
KG-A2C	10.8%	
DRRN	13.0%	

Example *Zork I* walkthrough state

- **Context:** [CLS] with a great effort, the rug is moved to one side of the room, revealing the dusty cover of a closed trap door. [SEP] open trapdoor [SEP] the door reluctantly opens to reveal a rickety staircase descending into darkness. you are carrying: a sword a nasty knife a rope a brass lantern a clove of garlic a jewel-encrusted egg living room you are in the living room. there is a doorway to the east, a wooden door with strange gothic lettering to the west, which appears to be nailed shut, a trophy case, and a rug lying beside an open trap door. [SEP]
- **GPT-2:** [east, west, down, up, north, south, open trophy case, wait, knock on door, take rug, southeast, enter trapdoor, out, drop sword, take rope, in, southwest, northwest, get rope, open case, get rug, search rug, enter trap, climb rope, northeast, take sword, move rug, take all, put sword in trapdoor, close trapdoor]
- n-gram: [north, east, south, west, up, down, open door, examine door, take all, unlock door, get all, close door, drop all, put all, tie rope, examine knife, take knife, examine case, examine sword, open case, examine rope, examine west, take rope, take sword, examine lantern, put knife, pull rope, take lantern, examine egg, put sword, get sword, put egg]
- Admissible actions (generated by handicap enumeration): [east, open egg with lantern, throw rope at egg, throw egg at knife, throw sword at egg, throw garlic at egg, throw lantern at egg, throw knife at egg, throw knife at egg, throw lantern, put down all, put down rope, put down egg, put down sword, put down garlic, put down lantern, put down knife, close trap, take on egg, open case, turn on lantern, down]
- Walkthrough action: down

CALM Analysis: Admissibility of CALM Actions

- Evaluate admissibility of CALM actions on Jericho walkthroughs
- Precision and recall of admissible actions among top-k CALM actions (k=1...40)

k: tradeoff action diversity v. quality!
 After k=20, GPT-2 captures some complex/uncommon actions n-gram cannot



CALM (GPT-2) Analysis: Value of Data

- Two sources of data:
 - GPT-2 pretraining
 - ClubFloyd training
- 20%/50%: less ClubFloyd data
- -PT: GPT-2 randomly initialized
- +Jericho: add back 8 Jericho game scripts in ClubFloyd training



CALM (GPT-2) Analysis: RL

- No RL: cannot work
- Number of actions k:
 - k=30 is default
 - k=20/40: overall similar with k=30, different scores across games
 - k=10: significantly worse, ~ CALM (n- d gram)



Summary

- A language model approach to text game action generation
- RL perspective: action space reduction via LM pre-training
- Language perspective: generalized functional use across games
- Code and dataset: <u>https://github.com/princeton-nlp/calm-textgame</u>