

# Assessing and Enhancing Language Models for Complex and Robust Logical Reasoning over Natural Language

Speaker: Qiming Bao

Personal Website: <https://14h034160212.github.io/>

School: Strong AI Lab, NAOInstitute, University of Auckland

13 August 2025

DAAD AINeT fellow 2025 on Natural Language Processing Presentation



# Outline

- Motivation
- Research Gap
- Contributions
- Future Work

# Motivation

- Existing language models are challenged to effectively perform **complex logical reasoning in natural language**, particularly when confronted with **unbalanced distributions of reasoning depths** in multi-step and more real-world logical reasoning datasets. (IJCLR-NeSy 2022)
- One main reason existing language models struggle with complex natural language reasoning is the **lack of real-world, complex natural language reasoning datasets**, and it is challenging to obtain reliable data from the web for building expansive training datasets. (The Findings of ACL 2024)
- Furthermore, when large language models come out, they demonstrate evident improvement on the public logical reasoning datasets like ReClor, LogiQA and LogiQAv2, but whether this means those large language models have **strong and robust logical reasoning ability** remains to be seen. (LLM@IJCAI 2023)

# Multi-Step Deductive Reasoning Over Natural Language: An Empirical Study on Out-of-Distribution Generalisation

Authored by: Qiming Bao<sup>1</sup>, Alex Yuxuan Peng<sup>1</sup>, Tim Hartill<sup>1</sup>, Neşet Özkan Tan<sup>1</sup>, Zhenyun Deng<sup>1</sup>, Michael Witbrock<sup>1</sup>, Jiamou Liu<sup>1</sup>

<sup>1</sup>Strong AI Lab, NAOInstitute, Waipapa Taumata Rau - The University of Auckland

The 2nd International Joint Conference on Learning & Reasoning and 16th International Workshop on Neural-Symbolic Learning and Reasoning (IJCLR-NeSy 22)

<https://arxiv.org/abs/2207.14000>



# Outline

- Research Gap
  - Our Contribution
- Model Overview
- Experiment Result

# Research Gap

- Existing neural-symbolic methods such as Memory Attention Control networks (MAC) [3], Dynamic Memory Network [4], and DeepLogic [5] present visualized attention maps for multi-step logical reasoning tasks over symbolic logic programs. Can these methods be extended to natural language reasoning?
- **Our Contribution:** We utilized GloVe [2] to represent the input sequences, which consist of a context and a statement. The concatenated representations of the context and statement are fed into the model, such as a gated recurrent unit (GRU), for the RNN-based models. The model demonstrates strong performance, achieving over 85% test accuracy on first-order logic multi-step natural language reasoning tasks.

# Research Gap

- Although RNN-based models and pre-trained transformer models perform well in general, does that mean these models possess strong reasoning abilities?
- **Our Contribution:** We found that existing RNN-based baseline models and pre-trained transformer-based models do not perform well in terms of out-of-distribution (OOD) generalization on multi-step reasoning tasks under the following three scenarios:
  1. When the model is trained on data with shallow reasoning depths but tested on data with deeper reasoning depths.
  2. When the model is trained on synthetically generated data but tested on data paraphrased by humans.
  3. When the model is trained on unshuffled data but tested on shuffled data.

# Research Gap

- Existing multi-step deductive reasoning datasets, such as PARARULES and CONCEPTRULE V1 and V2, have unbalanced distributions over reasoning depths. Only a small portion of these datasets require deep reasoning ( $2 \leq \text{Depth} \leq 5$ ).
- **Our Contribution:** We present a larger deep multi-step logical reasoning dataset named PARARULE-Plus. This dataset specifically augments the data at deeper reasoning depths ( $2 \leq \text{Depth} \leq 5$ ). It has also been collected by OpenAI/Evals, and ChatGPT performs poorly on this dataset [2].



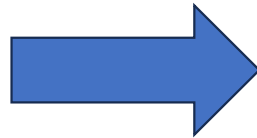
# From Symbolic Logic Program to Natural Language Reasoning

$p(X) : \neg q(X).$

$q(a).$

$?p(a).1$

$?p(b).0$



*(Input Facts:)* Alan is blue. Alan is rough. Alan is young.  
Bob is big. Bob is round.  
Charlie is big. Charlie is blue. Charlie is green.  
Dave is green. Dave is rough.

*(Input Rules:)* Big people are rough.  
If someone is young and round then they are kind.  
If someone is round and big then they are blue.  
All rough people are green.

Q1: Bob is green. True/false? [**Answer: T**]

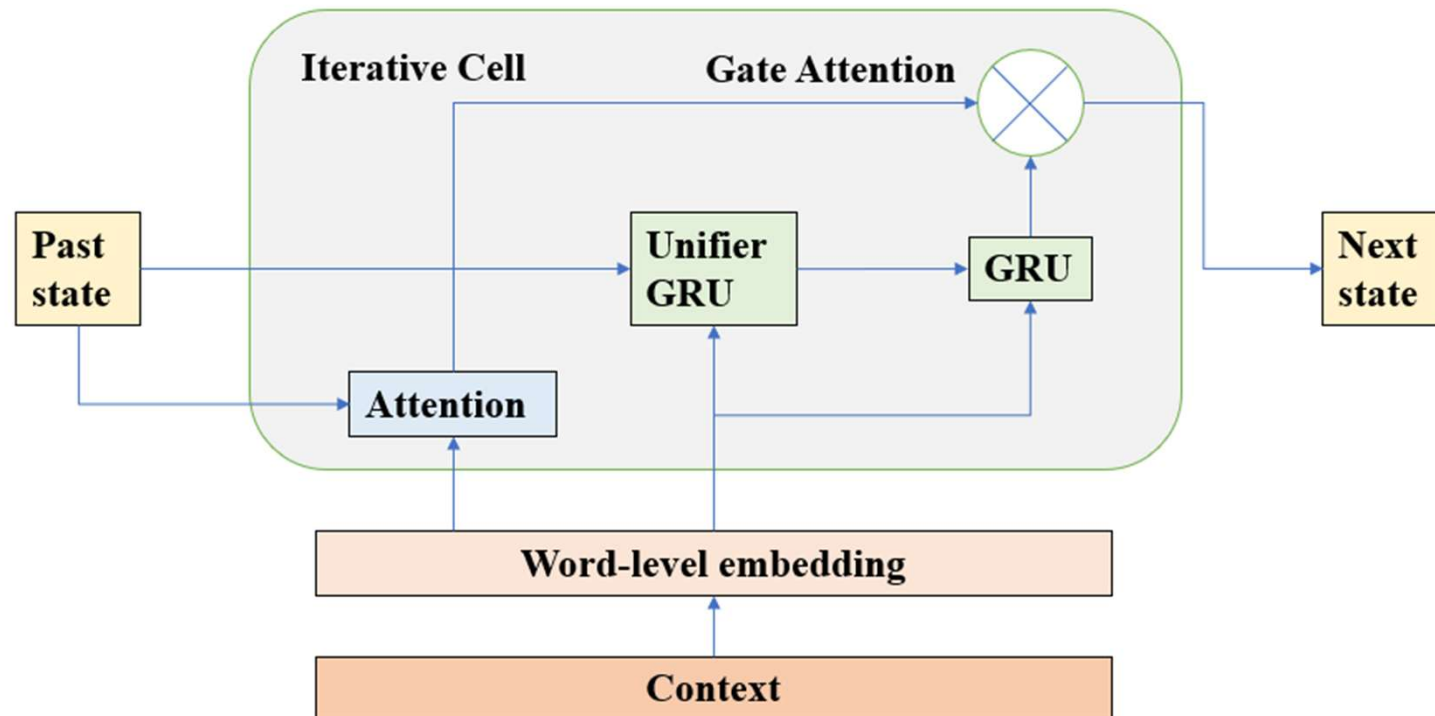
Q2: Bob is kind. True/false? [**F**]

Q3: Dave is blue. True/false? [**F**]



THE UNIVERSITY OF  
**AUCKLAND**  
Te Whare Wananga o Tāmaki Makaurau  
NEW ZEALAND

# Model Overview



# Experiment Result

**Table 4**

We use GloVe [16] as the word vector representation. We use PARARULES with all depths as the training set for all models and then test them on examples with different reasoning depths (D). Comparison among our IMA-GloVe-GA, IMA-GloVe, MAC-GloVe, DMN-GloVe, IMASM-GloVe, LSTM-GloVe, and RoBERTa-Large on PARARULES test sets with different reasoning depths.

Train ↓; Test →	D=1	D=2	D=3	$D \leq 3$	$D \leq 3 + \text{NatLang}$	$D \leq 5$	$D \leq 5 + \text{NatLang}$
IMA-GloVe	0.861	0.853	0.830	0.842	0.810	0.792	0.705
MAC-GloVe	0.792	0.776	0.750	0.763	0.737	0.701	0.652
DMN-GloVe	0.846	0.843	0.817	0.827	0.789	0.779	0.666
IMASM-GloVe	0.864	0.855	0.824	0.838	0.801	0.782	0.608
LSTM-GloVe	0.500	0.500	0.500	0.499	0.499	0.500	0.500
IMA-GloVe-GA	<b>0.950</b>	<b>0.943</b>	<b>0.919</b>	<b>0.927</b>	<b>0.883</b>	<b>0.879</b>	<b>0.741</b>
RoBERTa-Large	<b>0.986</b>	<b>0.985</b>	<b>0.977</b>	<b>0.979</b>	<b>0.972</b>	<b>0.967</b>	<b>0.949</b>

# Experiment Result

**Table 5**

IMA-GloVe, IMA-GloVe-GA, and RoBERTa-Large trained on CONCEPTRULES V1 (simplified / full) and tested on different test sets. Rules in CONCEPTRULES V1 Simplified are not shuffled, while CONCEPTRULES V1 full contains randomly shuffled rules. CONCEPTRULES V1 full has larger number of relations and entities than CONCEPTRULES V1 simplified.

Model	Train set	Test accuracy (Simplified Test set)	Test accuracy (Full Test set)
IMA-GloVe	Simplified	0.994	0.729
	Full	0.844	<b>0.997</b>
IMA-GloVe-GA	Simplified	<b>0.998</b>	<b>0.747</b>
	Full	0.851	<b>0.999</b>
RoBERTa-Large	Simplified	<b>0.997</b>	0.503
	Full	<b>0.927</b>	0.995

# Experiment Result

**Table 6**

IMA-GloVe, IMA-GloVe-GA, and RoBERTa-Large trained on CONCEPTRULES V2 (full) and tested on test sets that require different depths of reasoning.

Model	Test set	Mod1 Depth=1	Mod2 Depth=2	Mod3 Depth=3	Mod01 Depth $\leq$ 1	Mod012 Depth $\leq$ 2	Mod0123 Depth $\leq$ 3
IMA-GloVe	Depth=1	<b>0.999</b>	<b>0.998</b>	<b>0.990</b>	<b>0.997</b>	<b>0.998</b>	<b>0.997</b>
	Depth=2	<b>0.998</b>	<b>0.999</b>	<b>0.988</b>	<b>0.995</b>	<b>0.998</b>	<b>0.997</b>
	Depth=3	<b>0.997</b>	0.998	0.981	<b>0.991</b>	0.996	<b>0.997</b>
IMA-GloVe-GA	Depth=1	0.993	0.996	0.987	0.987	0.991	<b>0.997</b>
	Depth=2	0.993	<b>0.999</b>	0.974	0.986	0.991	0.995
	Depth=3	0.988	<b>1</b>	<b>0.994</b>	0.989	<b>0.997</b>	0.994
RoBERTa-Large	Depth=1	0.998	0.975	0.831	0.995	0.975	0.971
	Depth=2	0.997	0.972	0.885	0.993	0.972	0.965
	Depth=3	0.987	0.951	0.984	0.988	0.951	0.936



# Experiment Result

**Table 7**

RoBERTa-Large trained on PARARULES with different reasoning depths and tested on test sets that require different depths of reasoning. A bold number indicates the highest accuracy in a test set.

Model	Test set	Mod012 (Depth $\leq$ 2)	Mod0123 (Depth $\leq$ 3)	Mod0123Nat (Depth $\leq$ 3+NatLang)	Mod012345 (Depth $\leq$ 5)
RoBERTa-Large	Depth=0	<b>0.971</b>	0.946	0.968	0.953
	Depth=1	<b>0.943</b>	0.907	0.933	0.909
	Depth=2	<b>0.933</b>	0.902	0.932	0.902
	Depth=3	0.562	0.902	<b>0.926</b>	0.907
	Depth=4	0.481	0.863	<b>0.904</b>	0.888
	Depth=5	0.452	0.856	0.916	<b>0.933</b>
	NatLang	0.573	0.579	<b>0.962</b>	0.594

# Experiment Result

**Table 8**

RoBERTa-Large is fine-tuned on examples with different depths from PARARULES and also the entire PARARULE-Plus(PPT), and then is evaluated on test sets that require different depths of reasoning. The yellow background indicates improvement on accuracy after adding our PARARULE-Plus in the training process.

Model	Test set	Mod012 (Depth $\leq$ 2+PPT)	Mod0123 (Depth $\leq$ 3+PPT)	Mod0123Nat (Depth $\leq$ 3+NatLang+PPT)	Mod012345 (Depth $\leq$ 5+PPT)
RoBERTa-Large	Depth=0	0.946	0.901	0.965	<b>0.963 (+0.010)</b>
	Depth=1	0.877	0.847	<b>0.937 (+0.004)</b>	0.881
	Depth=2	0.868	0.873	<b>0.927</b>	0.839
	Depth=3	<b>0.771 (+0.209)</b>	0.862	<b>0.904</b>	0.826
	Depth=4	<b>0.675 (+0.194)</b>	0.852	<b>0.897</b>	0.832
	Depth=5	<b>0.661 (+0.209)</b>	<b>0.888 (+0.032)</b>	<b>0.923 (+0.007)</b>	<b>0.934 (+0.001)</b>
	NatLang	0.557	<b>0.593 (+0.014)</b>	<b>0.970 (+0.008)</b>	<b>0.649 (+0.055)</b>

# Abstract Meaning Representation-Based Logic-Driven Data Augmentation for Logical Reasoning

Authored by: **Qiming Bao**<sup>1,2</sup>, **Alex Yuxuan Peng**<sup>1</sup>, **Zhenyun Deng**<sup>3</sup>, **Wanjuan Zhong**<sup>4</sup>, **Gaël Gendron**<sup>1</sup>, **Timothy Pistotti**<sup>1</sup>, **Neşet Tan**<sup>1</sup>, **Nathan Young**<sup>1</sup>, **Yang Chen**<sup>1</sup>, **Yonghua Zhu**<sup>1</sup>, **Paul Denny**<sup>5</sup>, **Michael Witbrock**<sup>1</sup>, **Jiamou Liu**<sup>1</sup>

<sup>1</sup>Strong AI Lab, NAOInstitute, Waipapa Taumata Rau - The University of Auckland

<sup>2</sup>Xtracta, New Zealand

<sup>3</sup>Department of Computer Science and Technology, University of Cambridge, The United Kingdom

<sup>4</sup>School of Computer Science and Engineering, Sun Yat-Sen University, China

<sup>5</sup>School of Computer Science, The University of Auckland, New Zealand

The Findings of ACL 2024

<https://aclanthology.org/2024.findings-acl.353/>





# Outline

- Research Gap
  - Our Contribution
- System Architecture
- Experiment Results

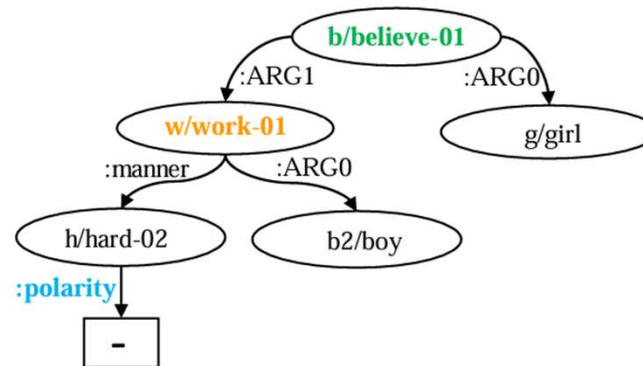
# Research Gap

- Enabling pre-trained large language models (LLMs) to reliably perform logical reasoning is an important step towards strong artificial intelligence [1]. The lack of available large real-world logical reasoning datasets means that LLMs are usually trained on more general corpora or smaller ones that do not generalise well.
- Logical reasoning is extremely important for solving problems in a robust, faithful and explainable way [2] [3], but because logical reasoning is complex for humans to understand and difficult to use for constructing data, there is exceptionally limited data. This implies that a scarcity of labeled datasets for logical reasoning persists in real-world scenarios. Consequently, it is not surprising that these pre-trained language models exhibit shortcomings in logical reasoning [4].

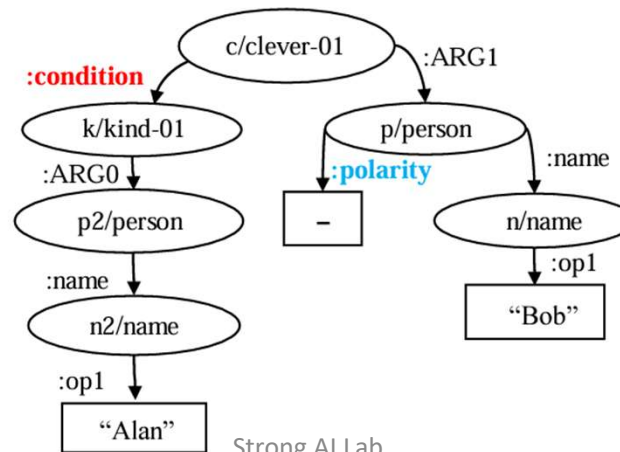
# Abstract Meaning Representation

S1: The girl **believes** that the boy **doesn't work** hard.

S2: That the boy **doesn't work** hard is what the girl **believes**.



S3: **If** Alan is kind, then Bob is **not** clever.



# Logical Reasoning Tasks

## Example Case

**Context:** If you have no keyboarding skills at all, you will not be able to use a computer. And if you are not able to use a computer, you will not be able to write your essays using a word processing program.

**Question:** If the statements above are true, which one of the following must be true?

### Options:

A. If you are not able to write your essays using a word processing program, you have no keyboarding skills.

*B. If you are able to write your essays using a word processing program, you have at least some keyboarding skills. ✓*

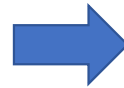
C. If you are not able to write your essays using a word processing program, you are not able to use a computer.

D. If you have some keyboarding skills, you will be able to write your essays using a word processing program.

$\alpha$  = you have keyboarding skills.

$\beta$  = you are able to use a computer.

$\gamma$  = you are able to write your essays using a word processing program.



Context:  $\neg \alpha \rightarrow \neg \beta$ ,  $\neg \beta \rightarrow \neg \gamma$

Option A:  $\neg \gamma \rightarrow \neg \alpha$

✓ Option B:  $\gamma \rightarrow \alpha + (\beta \rightarrow \alpha, \gamma \rightarrow \beta)$  using contraposition law

Option C:  $\neg \gamma \rightarrow \neg \beta$

Option D:  $\alpha \rightarrow \gamma$

A natural language logical reasoning reading comprehension example from ReClor[1].

Convert the natural language into logic symbols.

# Logical Equivalence Laws

## Definition 1: Contraposition law

$$(\mathcal{A} \rightarrow \mathcal{B}) \Leftrightarrow (\neg \mathcal{B} \rightarrow \neg \mathcal{A})$$

*If Alan is kind, then Bob is clever.  $\Leftrightarrow$  If Bob is not clever, then Alan is not kind.*

## Definition 2: Implication law

$$(\mathcal{A} \rightarrow \mathcal{B}) \Leftrightarrow (\neg \mathcal{A} \vee \mathcal{B})$$

*If Alan is kind, then Bob is clever.  $\Leftrightarrow$  Alan is not kind or Bob is clever.*

## Definition 3: Commutative law

$$(\mathcal{A} \wedge \mathcal{B}) \Leftrightarrow (\mathcal{B} \wedge \mathcal{A})$$

*Alan is kind and Bob is clever.  $\Leftrightarrow$  Bob is clever and Alan is kind.*

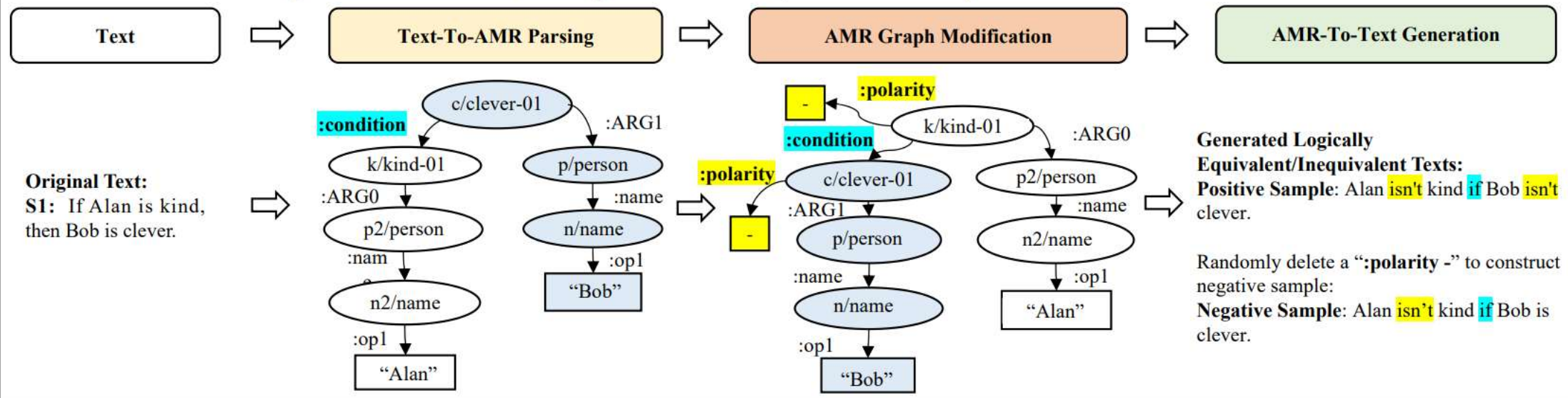
## Definition 4: Double negation law

$$\mathcal{A} \Leftrightarrow \neg \neg \mathcal{A}$$

*Alan is kind.  $\Leftrightarrow$  Alan is not unkind.*

# System Architecture

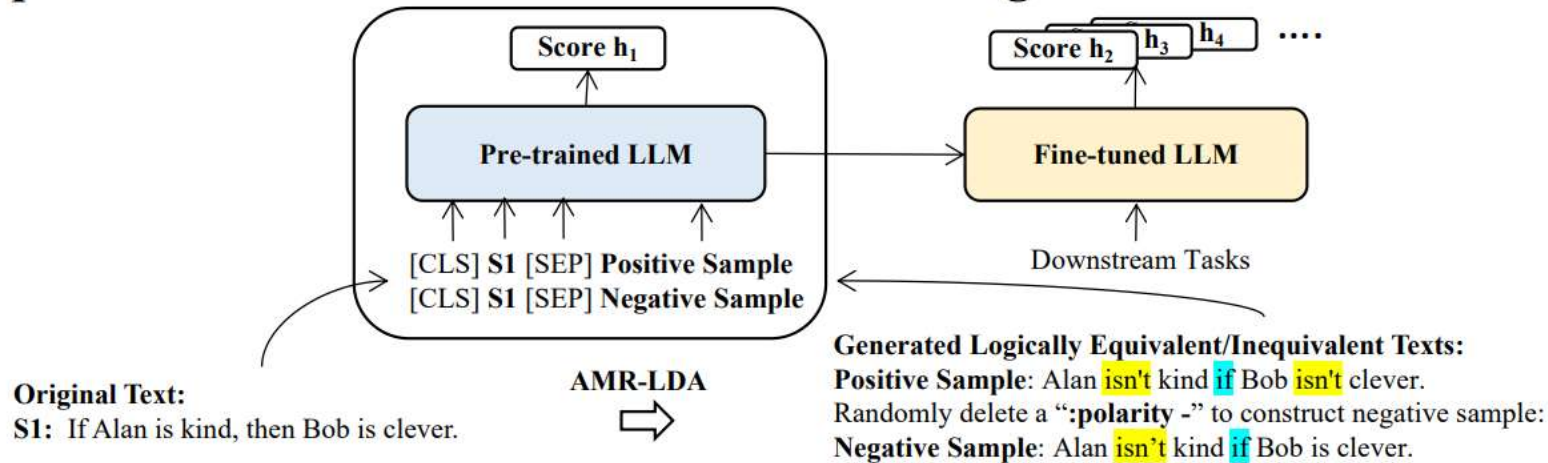
## 1. AMR-Based Logic-Driven Data Augmentation (AMR-LDA)





# System Architecture

## 2a. Logical-Equivalence-Identification Contrastive Learning for Discriminative LLM



# System Architecture

## 2b. Prompt Augmentation for Generative LLM

Context:  $\neg \alpha \rightarrow \neg \beta, \neg \beta \rightarrow \neg \gamma$

Option A:  $\neg \gamma \rightarrow \neg \alpha$

Option B:  $\gamma \rightarrow \alpha$

Option C:  $\neg \gamma \rightarrow \neg \beta$

Option D:  $\alpha \rightarrow \gamma$

AMR-LDA



Context:  $\neg \alpha \rightarrow \neg \beta, \neg \beta \rightarrow \neg \gamma$

Option A:  $\neg \gamma \rightarrow \neg \alpha$  + AMR-LDA extended option:  $\alpha \rightarrow \gamma$  + AMR-LDA extended context:  $\beta \rightarrow \alpha, \gamma \rightarrow \beta$

Option B:  $\gamma \rightarrow \alpha$  + AMR-LDA extended option:  $\neg \alpha \rightarrow \neg \gamma$  + AMR-LDA extended context:  $\beta \rightarrow \alpha, \gamma \rightarrow \beta$

Option C:  $\neg \gamma \rightarrow \neg \beta$  + AMR-LDA extended option:  $\beta \rightarrow \gamma$  + AMR-LDA extended context:  $\beta \rightarrow \alpha, \gamma \rightarrow \beta$

Option D:  $\alpha \rightarrow \gamma$  + AMR-LDA extended option:  $\neg \gamma \rightarrow \neg \alpha$  + AMR-LDA extended context:  $\beta \rightarrow \alpha, \gamma \rightarrow \beta$

$\alpha$  = you have keyboarding skills.

$\beta$  = you are able to use a computer.

$\gamma$  = you are able to write your essays using a word processing program.

Solution Path 1

Solution Path 2




Option B ✓



# Experiment Results

Models/ Datasets	ReClor				LogiQA		MNLI	MRPC	RTE	QNLI	QQP
	Dev	Test	Test-E	Test-H	Dev	Test			Eval		
RoBERTa	59.73	53.20	72.57	37.97	35.43	34.50	88.95	90.44	83.39	<b>94.73</b>	90.89
RoBERTa LReasoner-LDA	59.46	53.66	72.19	39.10	34.81	34.81	89.41	89.46	86.28	94.25	90.01
RoBERTa AMR-DA	58.66	53.93	66.81	<b>43.80</b>	36.45	37.22	89.74	90.44	86.28	94.42	92.06
RoBERTa AMR-LDA	<b>65.26</b>	<b>56.86</b>	<b>77.34</b>	40.77	<b>40.29</b>	<b>38.14</b>	<b>89.78</b>	<b>90.93</b>	<b>86.64</b>	94.49	<b>93.14</b>
DeBERTaV2	73.93	70.46	80.82	62.31	39.72	39.62	89.45	89.71	84.48	95.00	<b>92.54</b>
DeBERTaV2 LReasoner-LDA	75.73	70.70	84.08	60.17	30.87	28.51	89.23	89.95	87.00	95.15	92.50
DeBERTaV2 AMR-DA	79.06	75.90	84.62	69.04	29.95	30.10	<b>89.92</b>	89.71	83.39	95.02	92.42
DeBERTaV2 AMR-LDA	<b>79.40</b>	<b>77.63</b>	<b>85.75</b>	<b>71.24</b>	<b>42.34</b>	<b>39.88</b>	89.67	<b>90.20</b>	<b>88.09</b>	<b>95.24</b>	92.47

Table 2: Comparison between our proposed AMR-LDA and baseline models. We use RoBERTa-Large, DeBERTaV2-XXLarge as the pre-trained models. Our fine-tuned LLMs perform equally well or better than baseline methods.


**ReClor - A Reading Comprehension Dataset Requiring Logical Reasoning**
★ 40

Organized by: [ReClor Team](#)  
 Starts on: Jan 1, 2020 1:00:00 PM NZST (GMT + 13:00)  
 Ends on: Jan 1, 2100 12:59:59 PM NZST (GMT + 13:00)

Overview Evaluation Phases (A) Participate Leaderboard Discuss

Rank	Participant team	Test (t)	Test-E (t)	Test-H (t)	NA (t)	SA (t)	S (t)	W (t)	E (t)	I (t)	CMP (t)	MSS (t)	ER (t)	P (t)	D (t)	T (t)	R (t)	IF (t)
1	AMR-LDA Team	90.20	91.59	89.11	92.11	83.33	90.43	88.50	100.00	84.78	97.22	94.64	94.05	87.69	96.67	94.44	87.50	91
2	HFL & iFLYTEK (IDOL/Rational Reasoner)	80.60	87.73	75.00	86.84	90.00	84.04	72.57	76.92	58.70	86.11	73.21	82.14	76.92	80.00	86.11	81.25	83
3	MERIT (MERIT-deberta-v2-xxlarge)	79.30	85.23	74.64	85.09	83.33	82.98	71.68	76.92	65.22	83.33	73.21	76.19	80.00	80.00	88.89	78.13	81
4	Knowledge Model Team (Knowledge model)	79.20	91.82	69.29	89.47	80.00	76.60	68.14	92.31	63.04	94.44	78.57	78.57	78.46	76.67	97.22	84.38	76
5	AMR-LDA (DeBERTa-v2-xxlarge-AMR-LDA-Con)	77.20	86.14	70.18	83.33	76.67	79.79	68.14	84.62	52.17	88.89	80.36	75.00	75.38	80.00	88.89	71.88	78
6	LReasoner Team (LReasoner ensemble)	76.10	87.05	67.50	80.70	80.00	76.60	67.26	84.62	67.39	88.89	76.79	76.19	75.38	63.33	88.89	71.88	74

Models/Datasets	ReClor				LogiQA	
	Dev	Test	Test-E	Test-H	Dev	Test
GPT-3.5	57.02	56.20	59.31	<b>53.75</b>	37.63	37.32
+ CoT	34.80	25.80	27.50	24.46	23.96	24.57
+ AMR-DA	33.20	32.90	34.31	31.78	<b>40.55</b>	31.49
+ AMR-LDA	<b>58.62</b>	<b>56.69</b>	<b>60.90</b>	53.39	<b>40.55</b>	<b>39.47</b>
GPT-4	87.35	89.60	90.90	88.57	43.24	53.88
+ CoT	37.00	24.80	26.13	23.75	23.50	27.03
+ AMR-DA	85.00	85.60	86.36	85.00	51.30	56.06
+ AMR-LDA	<b>87.73</b>	<b>90.20</b>	<b>91.59</b>	<b>89.11</b>	<b>51.92</b>	<b>58.06</b>

Table 3: Comparison of Chain-of-Thought Prompting (CoT), AMR-DA, and AMR-LDA on GPT-3.5 and GPT-4, and between GPT-3.5 and GPT-4 alone, for evaluation on the ReClor and LogiQA test sets.

# Experiment Results

Models/Datasets	RoBERTa AMR-LDA	RoBERTa LReasoner-LDA
Depth=1	100.00	100.00
Depth=1 (with altered rules)	<b>100.00</b>	99.87
Depth=2	100.00	100.00
Depth=2 (with altered rules)	<b>99.73</b>	74.00

Table 4: Comparison between AMR-LDA and LReasoner-LDA with RoBERTa-Large on PARARULE-Plus and PARARULE-Plus (with altered rules). Depth=1 means that only one rule was used to infer the answer. Depth=1 (with altered rules) means one of the rules has been altered using logical equivalence law.

# Experiment Results

Models/Datasets	ReClor				LogiQA	
	Dev	Test	Test-E	Test-H	Dev	Test
DeBERTaV2-XXLarge	73.93	70.46	80.82	62.31	39.72	39.62
+ AMR-LDA-1:1	78.80	76.10	<b>84.77</b>	69.28	40.55	41.47
+ AMR-LDA-1:2	80.20	<b>76.40</b>	<b>84.77</b>	<b>69.82</b>	<b>47.00</b>	<b>43.93</b>
+ AMR-LDA-1:3	<b>81.20</b>	75.70	84.09	69.10	42.70	41.01
DeBERTaV2-XXLarge + MERIt-1:3	80.20	75.80	85.00	68.57	37.32	42.39
+ AMR-LDA-Con-1:3	<b>82.60</b>	76.60	86.13	<b>69.10</b>	<b>45.00</b>	43.01
+ AMR-LDA-Merged-1:3	81.80	<b>76.90</b>	<b>87.50</b>	68.57	44.54	<b>45.62</b>
DeBERTaV2-XXLarge + IDoL	77.60	74.50	82.95	67.85	39.78	40.24
+ AMR-LDA-Con-1:3	79.20	<b>77.00</b>	85.68	<b>70.17</b>	<b>47.61</b>	<b>44.54</b>
+ AMR-LDA-Merged-1:3	<b>79.40</b>	75.60	<b>86.36</b>	67.14	41.93	41.32

Table 6: An experiment to assess how positive:negative sample ratios affect downstream tasks. AMR-LDA 1:1 means the ratio of positive and negative samples is 1:1.

Models/Datasets	Con	Con-dou	Con-dou imp	Con-dou imp-com
<i>RoBERTa-Large as backbone model</i>				
ReClor	60.40	60.80	<b>61.80</b>	59.80
LogiQA	37.78	33.17	33.94	<b>38.70</b>
MNLI	89.55	<b>90.15</b>	89.68	89.78
MRPC	90.69	89.22	90.44	<b>90.93</b>
RTE	81.23	85.20	84.84	<b>86.64</b>
QNLI	94.16	94.05	<b>94.51</b>	94.49
QQP	92.12	89.88	92.06	<b>93.14</b>
<i>DeBERTaV2-XXLarge as backbone model</i>				
ReClor	<b>81.80</b>	72.20	79.40	78.80
LogiQA	32.25	<b>45.46</b>	38.24	40.55
<i>DeBERTa-Large as backbone model</i>				
MNLI	<b>90.80</b>	90.59	90.68	89.67
MRPC	<b>90.20</b>	88.48	89.95	<b>90.20</b>
RTE	84.84	87.36	85.56	<b>88.09</b>
QNLI	<b>95.28</b>	95.04	94.97	95.24
QQP	92.33	92.40	92.29	<b>92.47</b>

Table 5: An experiment to assess the influence of different logical equivalence laws on downstream logical reasoning and natural language inference tasks. “Con”, “dou”, “imp” and “com” are the abbreviation for contraposition law, double negation law, implication law and commutative law. “Con-dou” denotes data constructed using both the contraposition law and the double negation law. Other terms are derived in a similar manner.

# Useful Links



Project code



#1 on ReClor Leaderboard



Model Weights

Our AMR-LDA has been open-sourced in the project code, and the model weights have been released.

Welcome for more discussion and collaboration!

# A Systematic Evaluation of Large Language Models on Out-of-Distribution Logical Reasoning Tasks

Authored by: **Qiming Bao**<sup>1,2</sup>, **Gaël Gendron**<sup>1</sup>, **Alex Yuxuan Peng**<sup>1</sup>, **Wanjun Zhong**<sup>3</sup>, **Neşet Tan**<sup>1</sup>, **Yang Chen**<sup>1</sup>, **Michael Witbrock**<sup>1</sup>, **Jiamou Liu**<sup>1</sup>

<sup>1</sup>Strong AI Lab, NAOInstitute, Waipapa Taumata Rau - The University of Auckland

<sup>2</sup>Xtracta, New Zealand

<sup>3</sup>School of Computer Science and Engineering, Sun Yat-Sen University, China

The first edition of the Symposium on Advances and Open Problems in Large Language Models (**LLM@IJCAI'23**)

<https://arxiv.org/abs/2310.09430>

# Outline

- Research Gap
- System Architecture
- Experiment Results



# Research Gap

- Although large language models such as ChatGPT and GPT-4 perform well on public logical reasoning leaderboards like ReClor, does this mean these models possess strong logical reasoning abilities?
- **Our Contribution:**
- We find that existing large language models, like ChatGPT and GPT-4, perform well on the original publicly available logical reasoning datasets. However, their performance on our out-of-distribution test examples is poor, suggesting that the models might have seen these datasets during training and have failed to acquire generalized logical reasoning capabilities.
- We present a new set of challenging logical reasoning datasets, which involve shuffling the order of options and changing the correct option to "none of the other options are correct." These datasets have been collected by OpenAI/evals [1].

# Research Gap

- How to improve large language models' performance on complex logical reasoning tasks?
- **Our Contribution:**
- Relying solely on prompting techniques, such as chain-of-thought prompting, does not provide significant benefits for large language models when handling our newly presented, challenging logical reasoning tasks.
- We found that logic-driven data augmentation for instruction fine-tuning and prompting can help models improve their general performance on task structure variations in logical reasoning.



# OOD Logical Reasoning Tasks

## Example Case

**Context:** If you have no keyboarding skills at all, you will not be able to use a computer. And if you are not able to use a computer, you will not be able to write your essays using a word processing program.

**Question:** If the statements above are true, which one of the following must be true?

**Options:**

- A. If you are not able to write your essays using a word processing program, you have no keyboarding skills.
- B. If you are able to write your essays using a word processing program, you have at least some keyboarding skills. ✓*
- C. If you are not able to write your essays using a word processing program, you are not able to use a computer.
- D. If you have some keyboarding skills, you will be able to write your essays using a word processing program.

## 1. Shuffle Option Order

### Example Case

**Context:** If you have no keyboarding skills at all, you will not be able to use a computer. And if you are not able to use a computer, you will not be able to write your essays using a word processing program.

**Question:** If the statements above are true, which one of the following must be true?

**Options:**

- A. If you are not able to write your essays using a word processing program, you have no keyboarding skills.
- B. If you are not able to write your essays using a word processing program, you are not able to use a computer.
- C. If you are able to write your essays using a word processing program, you have at least some keyboarding skills. ✓*
- D. If you have some keyboarding skills, you will be able to write your essays using a word processing program.

## 2. Replace the correct option

### Example Case

**Context:** If you have no keyboarding skills at all, you will not be able to use a computer. And if you are not able to use a computer, you will not be able to write your essays using a word processing program.

**Question:** If the statements above are true, which one of the following must be true?

**Options:**

- A. If you are not able to write your essays using a word processing program, you have no keyboarding skills.
- B. None of the other options are correct. ✓*
- C. If you are not able to write your essays using a word processing program, you are not able to use a computer.
- D. If you have some keyboarding skills, you will be able to write your essays using a word processing program.

## 3. Shuffle and replace the correct option

### Example Case

**Context:** If you have no keyboarding skills at all, you will not be able to use a computer. And if you are not able to use a computer, you will not be able to write your essays using a word processing program.

**Question:** If the statements above are true, which one of the following must be true?

**Options:**

- A. If you are not able to write your essays using a word processing program, you have no keyboarding skills.
- B. If you are not able to write your essays using a word processing program, you are not able to use a computer.
- C. None of the other options are correct. ✓*
- D. If you have some keyboarding skills, you will be able to write your essays using a word processing program.



THE UNIVERSITY OF  
**AUCKLAND**  
Te Whare Wananga o Tāmaki Makaurau  
NEW ZEALAND

# System Architecture

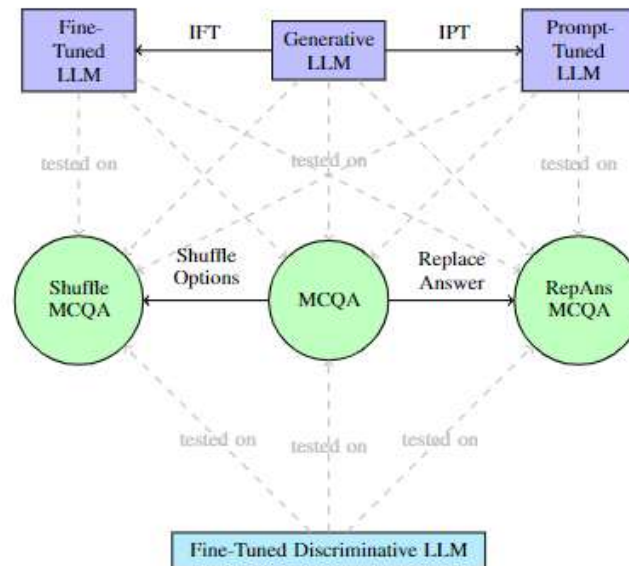


Figure 2: We perform instruction fine-tuning and prompting on the generative large language models and test them on several datasets. We also test fine-tuned discriminative large language models. We use Multiple-Choices Question Answering (MCQA) datasets and generate new distributions by shuffling the order of options and removing some answers. Square represent models, blue square represent generative models and cyan square represent classification engines, and green circles represent datasets.



# Experiment Results

Datasets → Models ↓	ReClor				LogiQA				LogiQAv2			
	Original	Shuffle Order	Replace Answer	Shuffle RepAns	Original	Shuffle Order	Replace Answer	Shuffle RepAns	Original	Shuffle Order	Replace Answer	Shuffle RepAns
<i>Zero-shot evaluation</i>												
Alpaca-7B	0.0020	0.0060	0.0100	0.0120	0.0122	0.0122	0.0107	0.0121	0.0216	0.0165	0.0095	0.0121
Vicuna-7B	0.0960	0.1120	0.0740	0.0640	0.2027	0.2135	0.1735	0.1784	0.0834	0.0618	0.0541	0.0121
GPT-3.5	0.5702	0.5734	0.1919	0.1847	0.3763	0.3946	0.2449	0.2583	0.5094	0.2695	0.2675	0.2583
GPT-4	0.8735	0.8405	0.1454	0.1312	0.4324	0.5283	0.1007	0.1686	0.5230	0.2616	0.1731	0.1686
<i>ReClor/LogiQA/LogiQAv2 single training set</i>												
Alpaca-7B-IFT	0.1680	0.5280	0.2360	0.2720	0.1105	0.3486	0.2841	0.2273	0.1912	0.2122	0.3658	0.1548
Vicuna-7B-IFT	0.3040	0.1760	0.0320	0.0420	0.2503	0.1689	0.0706	0.1198	0.1899	0.1746	0.1797	0.1784
LReasoner	0.7320	0.7100	0.2320	<b>0.3420</b>	0.4147	0.4316	<b>0.5176</b>	<b>0.5176</b>	0.5685	0.5685	<b>0.4263</b>	<b>0.4263</b>
MERIt	0.7960	0.7960	0.2580	0.2460	0.3794	0.3809	0.2657	0.2703	0.7144	0.7144	0.1873	0.1873
AMR-LE	0.8120	0.8120	<b>0.3360</b>	0.3360	0.4270	0.4301	0.1720	0.1720	0.6985	0.6978	0.1440	0.1440
<i>ReClor + LogiQA + LogiQAv2 merged training set</i>												
Alpaca-7B-IFT	0.7100	0.6560	0.1380	0.1140	0.6651	0.4854	0.2718	0.1351	0.6411	0.2160	0.1956	0.1128
Vicuna-7B-IFT	0.3900	0.4040	0.1500	0.1060	0.5453	0.3840	0.2273	0.1490	0.4913	0.1816	0.1708	0.1121
MERIt	0.9660	0.9660	0.2440	0.2440	0.7311	0.7342	0.2119	0.2119	0.8655	0.8661	0.2625	0.2625
AMR-LE	<b>0.9700</b>	<b>0.9700</b>	0.2900	0.2900	<b>0.7557</b>	<b>0.7588</b>	0.2549	0.2549	<b>0.8744</b>	<b>0.8744</b>	0.3212	0.3212

Table 4: Accuracy of large language models on logical reasoning tasks. The first block represents generative large language models tested in zero-shot settings. We compare them against models improved with instruction fine-tuning (IFT) on various training sets (separate training sets for the second block and merged training set for the third block). In the second block, models are fine-tuned on the original training dataset as they are evaluated on (e.g. fine-tuned on original ReClor training set and evaluated on ReClor validation set and other validation sets). In the third block, models are fine-tuned on a merged training set composed of all original training sets without our new datasets. Alpaca-7B and Vicuna-7B are trained using IFT fine-tuning and LReasoner, MERIt and AMR-LE are fine-tuned in the standard way.

# Experiment Results

Models	ReClor Shuffle RepAns	LogiQA Shuffle RepAns	LogiQAv2 Shuffle RepAns
<i>Zero-shot evaluation</i>			
Alpaca-7B	0.0120	0.0230	0.0121
Alpaca-7B-CoT	0.0120	0.0337	0.0152
Vicuna-7B	0.0640	0.1797	0.1784
Vicuna-7B-CoT	0.1320	0.1674	0.1593
GPT-3.5	<b>0.1847</b>	0.2286	<b>0.2583</b>
GPT-3.5-CoT	0.1088	0.1674	0.1722
GPT-4	0.1312	0.1626	0.1686
GPT-4-CoT	0.1816	<b>0.2523</b>	0.2177

Table 5: Comparison between base models and models prompted using Chain-of-Thought (CoT).

Models	ReClor Shuffle RepAns	LogiQA Shuffle RepAns	LogiQAv2 Shuffle RepAns
<i>Zero-shot evaluation</i>			
Alpaca-7B	0.0120	0.0121	0.0121
GPT-3.5	0.1847	0.2583	0.2583
GPT-4	0.1312	0.1686	0.1686
<i>ReClor/LogiQA/LogiQAv2 single training set</i>			
Alpaca-7B-IFT	0.2720	0.2273	0.1548
+ AMR-LE	0.0440	0.0522	0.0548
<i>ReClor + LogiQA + LogiQAv2 merged training set</i>			
Alpaca-7B-IFT	0.1140	0.1351	0.1128
+ AMR-LE	0.0060	0.0245	0.0197
<i>Prompt augmentation using AMR-LE</i>			
Alpaca-7B-IPT-LDA	0.0300	0.0368	0.0331
Alpaca-7B-IFT-LDA	0.4800	0.3686	0.2237
GPT-3.5-IPT-LDA	0.3667	0.4685	0.4971
GPT-4-IPT-LDA	<b>0.8766</b>	<b>0.5510</b>	<b>0.7027</b>

Table 6: Accuracy of evaluated models when adding AMR-LE’s logic-driven augmented data into the training set. We evaluate Alpaca-7B after instruction fine-tuning.



# Experiment Results

Datasets →												
Perturbation Ratio ↓	ReClor				LogiQA				LogiQAv2			
	Original	Shuffle Order	Replace Answer	Shuffle RepAns	Original	Shuffle Order	Replace Answer	Shuffle RepAns	Original	Shuffle Order	Replace Answer	Shuffle RepAns
<i>ReClor/LogiQA/LogiQAv2 single training set with different ratio of data perturbation (Shuffle-RepAns)</i>												
0%	0.1680	<b>0.5280</b>	0.2360	0.2720	0.1105	<b>0.3486</b>	0.2841	0.2273	0.1912	<b>0.2122</b>	0.3658	0.1548
5%	0.3340	0.3720	0.1560	0.1720	0.1490	0.1351	0.0998	0.0921	0.2695	0.1516	0.1338	0.1121
10%	<b>0.4140</b>	0.4320	0.2040	0.2380	<b>0.3072</b>	0.2826	0.2350	0.2442	0.2262	0.0956	0.1963	0.1727
15%	0.3620	0.3860	<b>0.3060</b>	<b>0.3340</b>	0.1904	0.2027	0.2795	0.2319	<b>0.3537</b>	0.1778	0.2001	0.1727
50%	0.1540	0.1400	0.1660	0.1640	0.0430	0.0537	<b>0.6728</b>	<b>0.6559</b>	<b>0.3537</b>	0.2096	<b>0.7686</b>	<b>0.7915</b>

Table 7: Accuracy of Alpaca-7B model for transfer learning scenarios and different perturbation ratio applied to the training set. To make a fair comparison, We ensure that the size of each training set is consistent.

Datasets	Train	Validation	Test
ReClor	4638	500	1000
LogiQA	7376	651	651
LogiQA-v2	12567	1569	1572

Table 5: Number of samples in the training, validation, and test set, for ReClor, LogiQA and LogiQA-v2.



Models	ReClor Shuffle RepAns	LogiQA Shuffle RepAns	LogiQAv2 Shuffle RepAns
<i>Zero-shot evaluation</i>			
LLaMA-7B	<b>0.1260</b>	0.1167	0.1128
LLaMA-13B	0.0660	0.1167	0.1013
LLaMA-30B	0.0360	0.1290	<b>0.1172</b>
LLaMA-65B	0.0720	<b>0.1397</b>	0.1159

Table 8: Comparison between multiple LLaMA model sizes on logical reasoning tasks with structure variations.

# Useful Links



Project Code



Strong AI Lab



LIU AI Lab



LR-MRC-Plus

Our proposed three logical reasoning reading comprehension datasets (ReClor-Plus, LogiQA-Plus and LogiQA-v2-plus have been collected by OpenAI Evals)

Welcome for more discussion and collaboration!



# Selected Publication List

- Qiming Bao, Alex Peng, Zhenyun Deng, Wanjun Zhong, Gaël Gendron, Neşet Tan, Nathan Young, Yang Chen, Yonghua Zhu, Michael Witbrock, Jiamou Liu. *Abstract Meaning Representation-Based Logic-Driven Data Augmentation for Logical Reasoning.*, the Findings of [ACL-24](#) [[#1 on the ReClor Leaderboard](#)] [[Paper link](#)] [[Source code](#)]
- Qiming Bao, Juho Leinonen, Alex Peng, Gaël Gendron, Wanjun Zhong, Tim Pistotti, Paul Denny, Michael Witbrock, Jiamou Liu. *Exploring Iterative Enhancement for Improving Learnersourced with Large Language Models*, The Proceeding of AACL 2025 ([AACL-EACL 2025](#)) [[Paper link](#)]
- Qiming Bao, Juho Leinonen, Alex Yuxuan Peng, Wanjun Zhong, Tim Pistotti, Alice Huang, Paul Denny, Michael Witbrock, Jiamou Liu. *Exploring Iterative Enhancement for Improving Learnersourced Multiple-Choice Question Explanations with Large Language Models*, [AGI@ICLR 2024](#) [[Paper link](#)] [[Source code](#)]
- Qiming Bao, Gaël Gendron, Alex Peng, Neşet Tan, Michael Witbrock, Jiamou Liu. *A Systematic Evaluation of Large Language Models on Out-of-Distribution Logical Reasoning Tasks.*, [LLM@IJCAI'23](#) [[Paper link](#)] [[Source code](#)]
- Qiming Bao, Alex Peng, Tim Hartill, Neşet Tan, Zhenyun Deng, Michael Witbrock, Jiamou Liu. *Multi-Step Deductive Reasoning Over Natural Language: An Empirical Study on Out-of-Distribution Generalisation*, [IJCLR-NeSy-22](#) [[Paper link](#)] [[Source code and dataset](#)] [[Presentation recording](#)]