
Adaptive Pre-training of Language Models for Better Logical Reasoning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Logical reasoning of text is an important ability that requires understanding the
2 logical information present in the text and reasoning through them to infer new
3 conclusions. Prior works on improving the logical reasoning ability of language
4 models require complex processing of training data (e.g., aligning symbolic knowl-
5 edge to text), yielding task-specific data augmentation solutions that restrict the
6 learning of general logical reasoning skills. In this work, we propose AERIE,
7 an adaptively pre-trained language model that has improved logical reasoning
8 abilities. We select a subset of Wikipedia, based on a set of logical inference key-
9 words, for continued pretraining of a language model. We use two self-supervised
10 loss functions: a modified masked language modeling loss where only specific
11 parts-of-speech words, that would likely require more reasoning than basic lan-
12 guage understanding, are masked, and a sentence classification loss that teaches
13 the model to distinguish between entailment and contradiction types of sentences.
14 The proposed training paradigm is both simple and generalizable across tasks.
15 We demonstrate the effectiveness of AERIE by comparing it with prior baselines
16 on two logical reasoning datasets. AERIE performs comparably on ReClor and
17 outperforms baselines on LogiQA.

18 1 Introduction

19 Logical reasoning is an important ability of humans that helps us in making rational decisions based
20 on some known information. Recently, logical reasoning of text has seen an increasing focus as it is a
21 fundamental skill required to solve any downstream task that requires machine reading [Yu et al.,
22 2020, Liu et al., 2021]. In these datasets, the model needs to understand a given context, reason
23 about a question, and then select the correct answer from a set of options. With the advent of large
24 pre-trained language models (PLMs) in NLP [Devlin et al., 2019, Radford et al., 2019, Raffel et al.,
25 2020], understanding and improving the logical reasoning abilities of these models has become even
26 more important as these are increasingly being used across a wide variety of real-world tasks.

27 There have been some recent works on improving the logical reasoning abilities of PLMs [Wang et al.,
28 2022, Ouyang et al., 2022, Jiao et al., 2022]. These works typically generate a dataset containing
29 symbolic structures such as logical graphs from text, logical contrast sets, etc., and then train the LM
30 using custom loss objectives to learn logical reasoning abilities. While the performance improvements
31 achieved by these methods are encouraging, the proposed solutions generally require complex data
32 processing to generate the additional structural information (graphs, contrast data, etc.) required for
33 logical reasoning. Further, the loss functions proposed in these works are very specifically designed
34 in accordance to their respective data augmentation technique, and widely differs from the typical
35 masked language modeling loss used for LM pretraining [Devlin et al., 2019]. These complex

36 processing steps usually require task-specific design choices, which are not necessarily learning
37 generalizable logical reasoning ability that is reusable across different task formats. Also, it is unclear
38 if these specific inductive biases are indeed essential for improving the logical reasoning abilities in
39 language models, or a simpler approach is sufficient.

40 Prior works [Gururangan et al., 2020] have shown that continued domain-adaptive pretraining of
41 PLMs lead to performance gains on downstream tasks. Inspired by this, we propose AERIE, a
42 continued pretraining-based approach to inject logical reasoning abilities in language models. To
43 gather a dataset that can teach logical reasoning, we use a set of keywords to select a subset of
44 the Wikipedia, such that every sentence in the subset contains at least one of the keywords. These
45 keywords are chosen such that the sentences containing the keywords are more likely to elicit
46 reasoning when filling out masked tokens. We note that in contrast to previous works [Gururangan
47 et al., 2020], our method only requires selecting sentences from Wikipedia, eliminating the need
48 for extra domain-specific corpus. Secondly, we restrict the type of tokens being masked from *any*
49 random token, to only specific types of tokens based on the parts-of-speech of the word. This choice
50 is again based on increasing the likelihood of using logical reasoning to predict the masked word.
51 Lastly, we add a sentence-level classification loss to predict if the reasoning in the sentence conveys
52 an entailment or a contradiction. This enables the model to understand the differences between these
53 two types of logical reasoning.

54 To test AERIE, we evaluate it on two downstream logical reasoning tasks: ReClor and LogiQA, and
55 compare it with other baselines. We achieve state-of-the-art performance on LogiQA and comparable
56 performance on ReClor. This demonstrates that our simple approach is generalizable to different
57 datasets and enables the PLM to learn logical reasoning abilities.

58 2 Problem Statement

59 In this work, we study the problem of using logical reasoning to solve the task of multiple choice
60 question answering based on a given context. Formally, for a given context C , question Q , and a
61 list of K candidate answers $A = \{A_1, \dots, A_K\}$, the task is to select the correct answer A_y , where
62 $y \in [1, K]$. Getting to the right answer typically requires reasoning logically through the context
63 and then selecting the best answer for the question. Evaluation of a model is based on the accuracy
64 metric.

65 3 Method

66 In this section, we describe the details of our proposed approach. In AERIE, we use a keyword-based
67 dataset selection strategy to collect a dataset of reasoning-related sentences called IMPLICATION
68 (§3.1) and then continue training a pretrained model checkpoint using two loss functions jointly
69 (§3.2). This model is then finetuned on the training dataset of each task separately.

70 3.1 Dataset Selection

71 PLMs are typically trained on the data from the internet which helps them in learning the language
72 model and then they are finetuned on specific downstream datasets to specialize on a task [Devlin
73 et al., 2019, Radford et al., 2018, Raffel et al., 2020]. We hypothesize that using a training data that
74 contains more reasoning related sentences, rather than generic internet data, should help in improving
75 the logical reasoning abilities of the PLM. Although creating such a dataset can be a challenging task
76 in itself, in AERIE, we explore a simple and intuitive way to curate a set of such sentences. First, we
77 select logical keywords that are generally encountered in sentences with some implication. Broadly,
78 we categorize these keywords into two types:

- 79 • **Positive implication (Entailment):** These keywords are typically present in sentences
80 where the reason entails the inference. We consider the following keywords in this category:
81 “therefore”, “accordingly”, “so”, “thus”, “consequently”, “hence”, “thence”, “and so”, “for
82 this reason”, “in consequence”, “on account of”, “on the ”grounds”, “since”, “therefrom”,
83 “thereupon”, “to that end”, “whence”, and “wherefore”.
- 84 • **Negative implication (Contradiction):** In this category, the keywords are usually present
85 in sentences where the reason contradicts the inference. Here, we consider the following

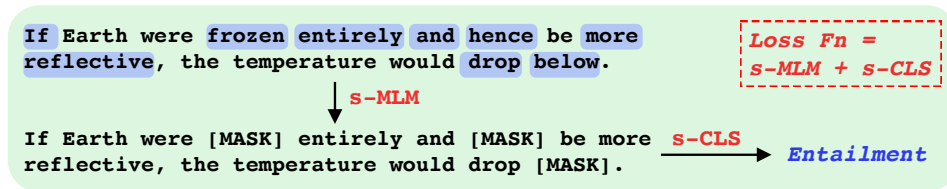


Figure 1: **Loss Functions in AERIE.** The s-MLM loss masks tokens from a specific set of POS tags (candidate tokens highlighted in blue), instead of any random token. The s-CLS loss classifies the masked sentence into one of two categories: entailment or contradiction. The overall loss function is the sum of both loss functions.

86 keywords: “but”, “although”, “however”, “nevertheless”, “on the other hand”, “still”,
87 “though”, and “yet”.

88 Next, we select sentences that contain at least one of the keywords. Specifically, we filter sentences
89 from Wikipedia ¹ such that they contain at least one of the keywords. We name this filtered version
90 of the Wikipedia as IMPLICATION. While this keyword-based filtering does not necessarily ensure
91 that the sentence has an implication statement, it increases the chances of such logically rich sentence
92 being present in the training set.

93 3.2 Loss Function Design

94 **Selective masked language modeling loss (s-MLM)** This is a modified version of the masked
95 language modeling (MLM) loss used in BERT [Devlin et al., 2019]. In the MLM loss, tokens in a
96 sentence are masked at random and the model learns to predict the masked tokens. While this helps
97 in learning a good language model, we hypothesize that not all masked tokens require similar degree
98 of reasoning to predict them. For example, most prepositions in a sentence are generally governed by
99 the English grammar. In contrast, some specific parts-of-speech (POS) tags such as adverbs require
100 more reasoning to predict the right token. Thus, in s-MLM, we mask out tokens that belong to a
101 specific set of POS tags. In AERIE, we mask tokens from the following POS tags [Honnibal and
102 Montani, 2017]: “ADJ”, “ADV”, “CONJ”, “CCONJ”, “PART”, “SCONJ”, and “VERB”.

103 **Sentence classification loss (s-CLS)** In addition to s-MLM, we add another auxiliary loss function
104 that predicts whether a sentence contains reasoning that entails or contradicts the inference. To predict
105 if a sentence is related to a positive or negative implications, a model would require strong logical
106 reasoning abilities. The labels for this loss is bootstrapped using the simple heuristic of whether the
107 specific type of keyword is present in the sentence. We note that although the keyword can be a direct
108 feature that can be used to predict the label, on average the keyword would be masked more often
109 due to our selective masking policy, leading to teaching the model some logical semantics.

110 4 Experimental Setup

111 Following prior works [Jiao et al., 2022], we evaluate AERIE on two logical reasoning datasets:
112 ReClor [Yu et al., 2020] and LogiQA [Liu et al., 2021]. Both the datasets are reading comprehension
113 style datasets, where the metric is the accuracy of the model in selecting the right answer for a
114 given context and question pair. We compare AERIE with three prominent baselines: LRReasoner
115 [Wang et al., 2022], Focal Reasoner [Ouyang et al., 2022], and MERIt [Jiao et al., 2022]. All
116 these baselines train a PLM using some additional data to improve logical reasoning abilities.

118 5 Results

119 **Overall Results** We use RoBERTa-Large pretrained check-
120 points as the starting point for AERIE and all the baselines.
121 In Table 1, we compare the performance of our method
122 with the baselines on the two logical reasoning datasets.
123 Overall, we observe that AERIE performs at par on ReClor
124 and outperforms all baselines on LogiQA.

| Model | ReClor | | LogiQA | |
|----------------|-------------|-------------|-------------|-------------|
| | Dev | Test | Dev | Test |
| RoBERTa | 62.6 | 55.6 | 35 | 35.3 |
| DAGN | 65.2 | 58.2 | 35.5 | 38.7 |
| DAGN (Aug) | 65.8 | 58.3 | 36.9 | 39.3 |
| LRReasoner | 64.7 | 62.4 | 38.1 | 40.6 |
| Focal Reasoner | 66.8 | 58.9 | 41.0 | 40.3 |
| MERIt | 66.8 | 59.6 | 40.0 | 38.9 |
| AERIE | 66.8 | 57.6 | 41.6 | 42.1 |

Table 1: Comparison of AERIE with other baselines on ReClor and LogiQA.

¹<https://huggingface.co/datasets/wikipedia>

125 **Ablation Studies** To study the effect of using IMPLICA-
 126 TION for continued pretraining along with the proposed loss functions, we first create RANDOM, a
 127 random subset of Wikipedia of similar size as that of IMPLICATION, and also consider using the
 128 standard masked language modeling (MLM) loss Devlin et al. [2019], where any token can be masked
 129 at random. The results of the ablation are shown in Table 2. We observe that using the IMPLICATION
 130 dataset leads to consistent improvements on both datasets, when compared to RANDOM dataset.
 131 Additionally, we find that both the s-MLM and s-CLS loss lead to improvements over MLM loss.
 132 Thus, this empirically justifies our choice of the dataset and loss functions proposed here.

133 6 Related Works

134 Reasoning in natural language
 135 has been a prevalent problem in
 136 NLP. In recent years, logical reason-
 137 ing in text has seen an increas-
 138 ing focus. ReClor [Yu et al.,
 139 2020] and LogiQA [Liu et al.,
 140 2021] are reading comprehension
 141 style datasets focused on ques-
 142 tions that require reasoning using
 143 information from a given context.

144 Wang et al. [2022] proposed LRReasoner, which parses symbolic logical structures from the train-
 145 ing data of ReClor for data augmentation using logical context extensions. Ouyang et al. [2022]
 146 constructed logical graphs using the chain of facts present in a task instance, and used GNNs to
 147 reason on the graph. Jiao et al. [2022] proposed MERIt, that used Wikipedia to generate sentence
 148 pairs for contrastive learning that are logically related, and trained the PLM using contrastive loss.
 149 Both LRReasoner and Focal Reasoner use data augmentation that are specific to the task being
 150 solved, making the pretraining process specific to the downstream dataset, and thus not generalizable
 151 across tasks. While MERIt addresses this issue by using Wikipedia to generate logical graphs, their
 152 contrastive loss formulation requires counterfactual data augmentation, that potentially distorts the
 153 factual knowledge present in the pretrained model. In contrast to prior works, we propose a simple
 154 continued pretraining strategy using minor modifications of standard masked language modeling
 155 loss [Devlin et al., 2019] and sentence classification loss to improve the logical reasoning ability of
 156 language models. Our approach is simple to integrate during pretraining, and is generalizable across
 157 tasks.

158 In a related line of work, a set of works [Clark et al., 2020, Saha et al., 2020, Tafjord et al., 2021,
 159 Sanyal et al., 2022b] used synthetically generated data to show that PLMs can perform complex
 160 deductive reasoning to predict entailment of a given hypothesis. While progress on these datasets are
 161 encouraging, some recent works have questioned if models are indeed robustly learning to perform
 162 logical reasoning Sanyal et al. [2022a].

163 7 Conclusion

164 In this paper, we proposed AERIE, an adaptive pre-trained language model with logical reasoning
 165 abilities. We use a subset of Wikipedia sentences for continued pretraining of the model using two
 166 self-supervised loss functions. The choice of the training dataset and loss functions are guided by
 167 the objective to include more reasoning related sentences and training signals, respectively. Through
 168 experiments on two logical reasoning datasets and ablation studies, we demonstrate the effectiveness
 169 of our proposed approach. Overall, we show that AERIE is a generalized solution to improving
 170 logical reasoning in language models.

171 References

172 Peter Clark, Oyvind Tafjord, and Kyle Richardson. Transformers as soft reasoners over language.
 173 In Christian Bessiere, editor, *Proceedings of the Twenty-Ninth International Joint Conference on*
 174 *Artificial Intelligence, IJCAI-20*, pages 3882–3890. International Joint Conferences on Artificial
 175 Intelligence Organization, 7 2020. Main track.

- 176 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep
177 bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of*
178 *the North American Chapter of the Association for Computational Linguistics: Human Language*
179 *Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota,
180 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1423. URL <https://aclanthology.org/N19-1423>.
- 182 Suchin Gururangan, Ana Marasović, Swabha Swayamdipta, Kyle Lo, Iz Beltagy, Doug Downey,
183 and Noah A. Smith. Don’t stop pretraining: Adapt language models to domains and tasks. In
184 *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages
185 8342–8360, Online, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.
186 acl-main.740. URL <https://aclanthology.org/2020.acl-main.740>.
- 187 Matthew Honnibal and Ines Montani. spaCy 2: Natural language understanding with Bloom
188 embeddings, convolutional neural networks and incremental parsing. To appear, 2017.
- 189 Fangkai Jiao, Yangyang Guo, Xuemeng Song, and Liqiang Nie. MERIt: Meta-Path Guided Con-
190 trastive Learning for Logical Reasoning. In *Findings of the Association for Computational*
191 *Linguistics: ACL 2022*, pages 3496–3509, Dublin, Ireland, May 2022. Association for Compu-
192 tational Linguistics. doi: 10.18653/v1/2022.findings-acl.276. URL <https://aclanthology.org/2022.findings-acl.276>.
- 194 Jian Liu, Leyang Cui, Hanmeng Liu, Dandan Huang, Yile Wang, and Yue Zhang. Logiqa: A
195 challenge dataset for machine reading comprehension with logical reasoning. In *Proceedings of*
196 *the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI’20*, 2021. ISBN
197 9780999241165.
- 198 Siru Ouyang, Zhuosheng Zhang, and hai zhao. Fact-driven logical reasoning, 2022. URL <https://openreview.net/forum?id=gKWxifgJVP>.
- 200 Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. Improving language
201 understanding by generative pre-training. URL [https://s3-us-west-2.amazonaws.com/openai-](https://s3-us-west-2.amazonaws.com/openai-assets/researchcovers/languageunsupervised/language_understanding_paper.pdf)
202 [assets/researchcovers/languageunsupervised/language_understanding_paper.pdf](https://s3-us-west-2.amazonaws.com/openai-assets/researchcovers/languageunsupervised/language_understanding_paper.pdf), 2018.
- 203 Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language
204 models are unsupervised multitask learners. 2019.
- 205 Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena,
206 Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified
207 text-to-text transformer. *Journal of Machine Learning Research*, 21(140):1–67, 2020. URL
208 <http://jmlr.org/papers/v21/20-074.html>.
- 209 Swarnadeep Saha, Sayan Ghosh, Shashank Srivastava, and Mohit Bansal. PProver: Proof generation
210 for interpretable reasoning over rules. In *Proceedings of the 2020 Conference on Empirical*
211 *Methods in Natural Language Processing (EMNLP)*, pages 122–136, Online, November 2020.
212 Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.9. URL <https://aclanthology.org/2020.emnlp-main.9>.
- 214 Soumya Sanyal, Zeyi Liao, and Xiang Ren. Robustlr: Evaluating robustness to logical perturbation
215 in deductive reasoning, 2022a. URL <https://arxiv.org/abs/2205.12598>.
- 216 Soumya Sanyal, Harman Singh, and Xiang Ren. FaiRR: Faithful and robust deductive reason-
217 ing over natural language. In *Proceedings of the 60th Annual Meeting of the Association for*
218 *Computational Linguistics (Volume 1: Long Papers)*, pages 1075–1093, Dublin, Ireland, May
219 2022b. Association for Computational Linguistics. doi: 10.18653/v1/2022.acl-long.77. URL
220 <https://aclanthology.org/2022.acl-long.77>.
- 221 Oyvind Tafjord, Bhavana Dalvi, and Peter Clark. ProofWriter: Generating implications, proofs,
222 and abductive statements over natural language. In *Findings of the Association for Compu-*
223 *tational Linguistics: ACL-IJCNLP 2021*, pages 3621–3634, Online, August 2021. Associa-
224 tion for Computational Linguistics. doi: 10.18653/v1/2021.findings-acl.317. URL <https://aclanthology.org/2021.findings-acl.317>.
- 225

- 226 Siyuan Wang, Wanjun Zhong, Duyu Tang, Zhongyu Wei, Zhihao Fan, Daxin Jiang, Ming Zhou,
227 and Nan Duan. Logic-driven context extension and data augmentation for logical reasoning of
228 text. In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 1619–1629,
229 Dublin, Ireland, May 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.
230 findings-acl.127. URL <https://aclanthology.org/2022.findings-acl.127>.
- 231 Weihao Yu, Zihang Jiang, Yanfei Dong, and Jiashi Feng. Reclor: A reading comprehension dataset
232 requiring logical reasoning. In *International Conference on Learning Representations (ICLR)*,
233 April 2020.