1  We thank the reviewers for their detailed comments and their useful suggestions. We are excited that **R1** and **R4** find
2  our work interesting, timely and novel, and that our results demonstrate the fundamental limitations of Transformer
3  Language Models (TLMs) reasoning abilities. We thank **R2** and **R3** for acknowledging that our experiments show
4  interesting results that open up questions for future research and spur potential modeling innovations. Some of the
5  concerns seen in the reviews were common to more than one reviewer, so we address these by topic below.

6  @**R1**, **R2**, **R3**, **R4** **Generalization to larger networks and different architectures**. We thank the reviewers for
7  suggesting to improve the experimental results by analyzing different model architectures, such as larger transformer
8  models (**R2**,**R3**), pre-trained language models (**R3**,**R4**), Transformer encoder-decoder (**R2**) and Graph Transformers
9  (**R1**). In this rebuttal, we report results on larger transformer models. We agree to provide results of pre-trained models
10 like GPT-2 and Transformer encoder-decoder models in the final submission, and if time permits Graph Transformers.
11 Regarding training using a larger Transformer model (**R2**,**R3**), we agree that 2.5M parameters is small compared to
12 more traditional transformer architectures such as GPT-X. To address this concern, we trained a 20 layer auto-regressive
13 network, resulting in 145M parameters. We observe that the generalization capacity of this network is similar (43%) to
14 the 2.5M parameter network trained on the same data (46%). Our preliminary investigation on pre-trained language
15 models (**R3**,**R4**), suggests that the pre-trained model has similar trends as training from scratch but we acknowledge
16 it needs further investigation. Due to specific limitations during training, preliminary investigation on Transformer
17 encoder-decoder (**R2**) suggests weaker generalization scores which we aim to investigate further.

18 @**R1**, **R2**, **R3**, **R4** **On the motivation for using TLMs**. This is an excellent question that we think will benefit the
19 understanding of all reviewers. We agree with **R2** that existing literature explores pre-training abilities of TLMs on large
20 natural language corpora. While training on massive data can give certain advantages with respect to understanding the
21 meanings of words, we conjecture that such data gives models much less experience with reasoning over long inference
22 chains. We study the less understood issues related to how well TLMs are able to perform long chains of reasoning.
23 Moreover, recent work such as LAMA, T5 and GPT3 suggest that language models can be treated as knowledge
24 bases. This directly motivates us to investigate if language models can also learn certain reasoning strategies. Studying
25 these abilities would enable future research in using these models as dynamic knowledge bases that could infer new
26 knowledge even when it is not "stored" directly (i.e. seen during pre-training). We will add this discussion to the paper.

27 @**R2**, **R3** **Natural Language results**. We thank you for highlighting the importance of results on natural language
28 stories. We acknowledge that generalization is weaker in this harder setting, and we confirmed that by performing
29 additional experiments on the natural language split. We still find that the proof resolution strategy influences the
30 generalization capacity of TLMs. In particular, the conclusion that models trained on long, exhaustive proofs generalize
31 better than short proofs still holds. We plan on moving discussion from the Appendix to the main section of the paper
32 along with additional results.

33 @**R2**, **R3** **On the complexity of the task**. We acknowledge that the dataset in use (CLUTRR) is a toy dataset. However,
34 this in turn allows us to carefully analyze and control the difficulty of the experiments. For instance, with 20 possible
35 entities ($k$ in Section 3.3) and 20 possible family kinship relationships, the model have to learn 8,000 possible triples.
36 Given that even in this simplistic setup the generalization performance is not positive, this warrants a deeper inspection
37 of reasoning mechanisms of TLMs. We plan on extending our experiments with other datasets in the future.

38 @**R2**, **R4** **On the issue of "all facts are seen" / "the correct answer is seen in the prefix"**. While our models have
39 seen all possible facts in all training proofs, the target answer to a (story, question) pair is not seen in the prefix given to
40 the model, unless the proof is explicitly given as in Section 4.3. Here, our experiments reveal that beyond 7-step proofs,
41 the copy mechanism learned by TLMs becomes unreliable due to positional token embeddings. Position-agnostic
42 embeddings could help in solving this issue, which we leave as an exercise for future work.

43 @**R4** **On the backward-chaining contradiction**. It is a great question why backward-chaining proofs are easier to
44 use in Section 4.1, but are harder to generate in Section 4.2. We also agree with **R4**, this is due to the fact that backward
45 chaining proofs contain the answer in the first proof step. Thus, there is a higher probability of the model to generate
46 this step correctly and then use it while predicting the answer. This explains why the answer accuracy of such model is
47 relatively high while their proof validity is low. We will note this phenomenon in the final version of the paper.

48 In general, we will fix typos and broken references (**R1**,**R2**,**R4**), clarify the presentation and some notations (**R1**,**R4**),
49 and expand on the background section by discussing theorem proving (**R1**,**R4**). We would like to thank again all the
50 reviewers for their time and effort in reading our paper and giving us good feedback and suggestions.