

Peer Review

Review of: "Metadata Conditioning Accelerates Language Model Pre-training"

Weiqiang Jin¹

1. Xi'an Jiaotong University, Xi'an, China

This paper introduces an innovative pre-training method called Metadata Conditioning then Cooldown (MeCo), which incorporates metadata into model training, significantly enhancing data efficiency while improving the controllability of language models. However, I still have some doubts and suggestions regarding the article, and I hope the author can improve the paper based on my thoughts:

1. The technical details in the paper, such as model architecture and training process, are relatively complex. To make these concepts more accessible to a wider audience, some expressions can be simplified. For example, instead of saying "We employ standard optimization settings for language models, i.e., AdamW optimizer and cosine learning rate schedule," it could be simplified to "We use standard optimization settings, such as the AdamW optimizer and cosine learning rate schedule."
2. Some sections of the paper have a slightly rigid structure, particularly when describing methods and experimental results. More transitional sentences could be added to improve the flow between paragraphs. For instance, after "MeCo substantially accelerates pre-training," a sentence like "This advantage is clearly observed across models of different scales" could help the flow feel more natural.
3. Although the experimental results in the paper are thoroughly analyzed, the lack of real-world application scenarios may make it harder for readers to connect with the findings. Adding examples of practical applications would make the paper more engaging.
4. The experimental results and tables are currently quite academic and straightforward. To make the content more approachable, some terms or complex descriptions could be rephrased with more intuitive metaphors.

5. The paper contains many figures and tables that provide detailed data, but summarizing the key points of each figure concisely could help readers quickly grasp the important information. For example, “The figure shows that, under the same training conditions, MeCo significantly improves the task performance of the 1.6B model” would help readers capture the key takeaway more quickly.

In conclusion, the MeCo method demonstrates how metadata can be cleverly used to accelerate pre-training and improve model performance without adding computational burden. This paper not only provides a practical solution for optimizing language models but also opens up new possibilities for model controllability. With its simple yet effective strategy, MeCo is a valuable work worth learning from.

Declarations

Potential competing interests: No potential competing interests to declare.