



Instruction Tuning With Loss Over Instructions

Zhengxiang Shi¹, Adam X. Yang², Bin Wu¹, Laurence Aitchison², Emine Yimaz¹, Aldo Lipani¹

¹University College London, ²University of Bristol

Takeaways

- If the number of training data is limited and completions are short, including the prompt loss during instruction tuning might be advantageous on various NLP and open-ended generation tasks.
- We identify two scenarios where including the prompt loss is particularly useful: (1) The ratio between instruction length and output length in the training data is high; and (2) The number of training examples is limited.
- The improvement stems from reducing the tendency to overfit, particularly under limited training resource conditions: Instruction tuning on brief outputs or a small amount of data can potentially lead to rapid overfitting.

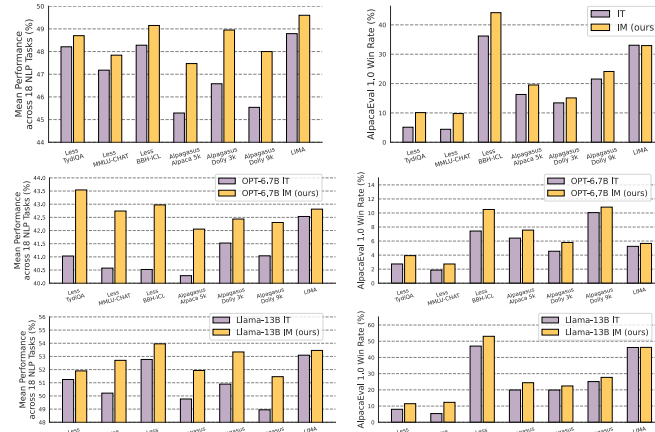
Limitations

- The success of our approach relies on the quality and diversity of the instructions and prompts in the training datasets.
- It is crucial to ensure that the instructions are ethically sound and free from harmful or biased content. Training on inappropriate or toxic instructions may result in undesirable outputs.

Main Experiments

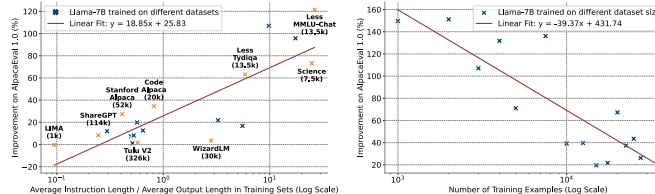
Main experiment 1: Performance differences between INSTRUCTION TUNING (IT, without the prompt loss) and INSTRUCTION MODELLING (IM, with the prompt loss) on 7 datasets.

Our findings: In many scenarios, IM can effectively improve the model performance on both NLP tasks (e.g., MMLU, TruthfulQA, and HumanEval) and open-ended generation benchmarks (e.g., MT-Bench and AlpacaEval).



Main experiment 2: Performance improvement, achieved by our approach INSTRUCTION MODELLING (IM) compared to INSTRUCTION TUNING (IT) on the AlpacaEval 1.0.

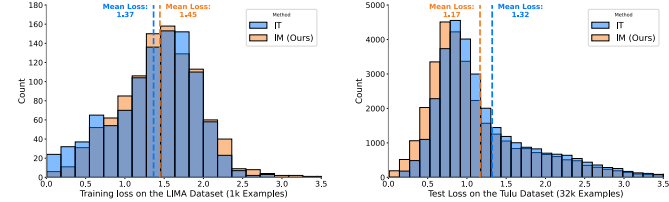
Our findings: We identify two key factors influencing the effectiveness of IM: (1) The ratio between instruction length and output length in the training data; and (2) The number of training examples.



Additional Experiments

Additional Experiment 1: Train and test loss analysis.

Our findings: IM has a higher train loss with lower test loss, suggesting that IM effectively mitigates the overfitting issues compared to IT.



Additional Experiment 2: Average BLEU Score comparison.

Our findings: IM produces outputs with less overlap with the ground truth outputs in training examples, indicating less overfitting.

	LIMA	Less Tydiqa	Less MMLU Chat	Less BBH ICL	AlpacaEval Alpaca 5k	AlpacaEval Dolly 9k	AlpacaEval Dolly 3k
IT	18.15	69.21	72.43	60.96	72.26	61.76	60.99
IM (ours)	17.30 _{0.85}	65.63 _{3.58}	69.20 _{3.23}	53.94 _{7.02}	70.50 _{1.76}	60.61 _{1.15}	59.04 _{1.95}

Additional Experiment 3: Can KL divergence loss, as regularization, easily address overfitting?

Our findings: (1) Incorporating KL Loss reduces overfitting and reduces the performance degradation on traditional NLP tasks; (2) KL Loss detrimentally affects model performance on open-ended generation tasks.

NLP Tasks	LIMA (1k)		ALPACAS DOLLY (9k)		
	LLAMA-2-7B-BASE	IT w/o KL Loss	IT w/ KL Loss	IT w/o KL Loss	IT w/ KL Loss
AlpacaEval 2.0	49.32	48.79 _{0.53}	49.26 _{0.06}	45.54 _{3.78}	49.31 _{0.01}
	0.01	2.58 _{72.57}	0.06 _{70.05}	2.28 _{72.27}	0.04 _{70.03}

Additional Experiment 4: Performance comparison of IM and IM +NEFTUNE on AlpacaEval 1.0 and various NLP benchmarks.

Our findings: Our proposed method IM could further improve the model performance with NEFTUNE.

	LIMA	Less Tydiqa	Less MMLU Chat	Less BBH ICL	AlpacaEval Alpaca 5k	AlpacaEval Dolly 9k	AlpacaEval Dolly 3k
	Mean Performance Across 18 NLP Tasks						
IM	32.94	10.10	9.78	44.15	19.52	30.77	15.11
IM +NEFTUNE	30.77 _{2.17}	23.41 _{713.31}	12.45 _{72.67}	48.25 _{74.10}	32.07 _{712.55}	38.28 _{77.51}	23.35 _{78.24}
IM	49.60	48.70	47.84	49.15	47.47	48.00	48.95
IM +NEFTUNE	49.47 _{10.13}	49.44 _{70.74}	47.73 _{10.11}	48.62 _{10.33}	48.70 _{71.23}	48.63 _{70.63}	49.54 _{70.59}