

NAACL 2025

The 5th Workshop on Insights from Negative Results in NLP

Proceedings of the Workshop

May 4, 2025

©2025 Association for Computational Linguistics

Order copies of this and other ACL proceedings from:

Association for Computational Linguistics (ACL)
317 Sidney Baker St. S
Suite 400 - 134
Kerrville, TX 78028
USA
Tel: +1-855-225-1962
acl@aclweb.org

ISBN 979-8-89176-240-4

Introduction

Publication of negative results is difficult in most fields, and the current focus on benchmark-driven performance improvement exacerbates this situation and implicitly discourages hypothesis-driven research. As a result, the development of NLP models often devolves into a product of tinkering and tweaking, rather than science. Furthermore, it increases the time, effort, and carbon emissions spent on developing and tuning models, as the researchers have little opportunity to learn from what has already been tried and failed.

The mission of the workshop on Insights from Negative Results in NLP is to provide a venue for many kinds of negative results, with the hope that they could yield useful insights and provide a much-needed reality check on the successes of deep learning models in NLP. In particular, we solicit the following types of contributions:

- broadly applicable recommendations for training/fine-tuning, especially if X that didn't work is something that many practitioners would think reasonable to try, and if the demonstration of X's failure is accompanied by some explanation/hypothesis;
- ablation studies of components in previously proposed models, showing that their contributions are different from what was initially reported;
- datasets or probing tasks showing that previous approaches do not generalize to other domains or language phenomena;
- trivial baselines that work suspiciously well for a given task/dataset;
- cross-lingual studies showing that a technique X is only successful for a certain language or language family;
- experiments on (in)stability of the previously published results due to hardware, random initializations, preprocessing pipeline components, etc;
- theoretical arguments and/or proofs for why X should not be expected to work;
- demonstration of issues with under-reporting of training details of pre-trained models, including test data contamination and invalid comparisons.

The fifth iteration of the Workshop on Insights from Negative Results attracted 23 submissions and 2 from ACL Rolling Reviews. We accepted 16 papers, resulting in 64% acceptance rate. We hope the workshop will continue to contribute to the many reality-check discussions on progress in NLP. If we do not talk about things that do not work, it is harder to see what the biggest problems are and where the community effort is the most needed

Table of Contents

<i>Challenging Assumptions in Learning Generic Text Style Embeddings</i>	
Phil Ostheimer, Marius Kloft and Sophie Fellenz	1
<i>In-Context Learning on a Budget: A Case Study in Token Classification</i>	
Uri Berger, Tal Baumel and Gabriel Stanovsky	7
<i>Reassessing Graph Linearization for Sequence-to-sequence AMR Parsing: On the Advantages and Limitations of Triple-Based</i>	
Jeongwoo Kang, Maximin Coavoux, Didier Schwab and Cédric Lopez	15
<i>Corrective In-Context Learning: Evaluating Self-Correction in Large Language Models</i>	
Mario S a n z - G u e r r e r o and Katharina Von Der Wense	24
<i>Do Prevalent Bias Metrics Capture Allocational Harms from LLMs?</i>	
Hannah Cyberek, Yangfeng Ji and David Evans	34
<i>Language-Specific Neurons Do Not Facilitate Cross-Lingual Transfer</i>	
Soumen Kumar Mondal, Sayambhu Sen, Abhishek Singhania and Preethi Jyothi	46
<i>Monte Carlo Sampling for Analyzing In-Context Examples</i>	
Stephanie Schoch and Yangfeng Ji	63
<i>Does Training on Synthetic Data Make Models Less Robust?</i>	
Lingze Zhang and Ellie Pavlick	79
<i>Bridging the Faithfulness Gap in Prototypical Models</i>	
Andrew Koulogorge, Sean Xie, Saeed Hassanpour and Soroush Vosoughi	86
<i>Aligning Sizes of Intermediate Layers by LoRA Adapter for Knowledge Distillation</i>	
Takeshi Suzuki, Hiroaki Yamada and Takenobu Tokunaga	100
<i>LLMs are not Zero-Shot Reasoners for Biomedical Information Extraction</i>	
Aishik Nagar, Viktor Schlegel, T h a n h - T u n g Nguyen, Hao Li, Yuping Wu, Kuluhan Binici and Stefan Winkler	106
<i>Exploring Limitations of LLM Capabilities with Multi-Problem Evaluation</i>	
Zhengxiang Wang, Jordan Kodner and Owen Rambow	121
<i>Exploring Multimodal Language Models for Sustainability Disclosure Extraction: A Comparative Study</i>	
Tanay Gupta, Tushar Goel and Ishan Verma	141
<i>Self Knowledge-Tracing for Tool Use (SKT-Tool): Helping LLM Agents Understand Their Capabilities in Tool Use</i>	
Joshua Vigel, Renpei Cai, Eleanor Chen, Anish Neema, Austen Liao, Kevin Zhu and Sean O'brien	150
<i>Error Reflection Prompting: Can Large Language Models Successfully Understand Errors?</i>	
Jason Li, Lauren Yraola, Kevin Zhu and Sean O'brien	157
<i>Evaluating Robustness of LLMs to Numerical Variations in Mathematical Reasoning</i>	
Yuli Yang, Hiroaki Yamada and Takenobu Tokunaga	171

Program

Sunday, May 4, 2025

09:00 - 09:10	<i>Opening Remarks</i>
09:10 - 09:50	<i>Technical session 1</i>
09:50 - 10:30	<i>Technical session 2</i>
10:30 - 11:30	<i>Coffee Break</i>
11:30 - 12:00	<i>Invited Talk 1</i>
12:30 - 14:00	<i>Lunch</i>
14:10 - 14:50	<i>Technical session 3</i>
14:50 - 15:30	<i>Technical session 4</i>
15:30 - 16:00	<i>Coffee Break</i>
16:00 - 16:30	<i>Invited Talk 2</i>
16:30 - 17:30	<i>Poster Session</i>