

[Re] A Reproduction of Automatic Multi-Label Prompting: Simple and Interpretable Few-Shot Classification

Victor Livernoche^{1,2, ID} and Vidya Sujaya^{1, ID}

¹McGill University, Montréal, Qc, Canada – ²Mila - Québec AI Institute, Montréal, Qc, Canada

Edited by

Koustuv Sinha,
Maurits Bleeker,
Samarth Bhargav

Received

04 February 2023

Published

20 July 2023

DOI

10.5281/zenodo.8173735

Reproducibility Summary

Scope of Reproducibility – We reproduce the results in the paper Automatic Multi-Label Prompting: Simple and Interpretable Few-Shot Classification by Wang et al. [1]. Using human prompt engineering can be long and expensive for text classification. In their paper they proposed an approach, called AMuLaP to do automatic label prompting for few-shot classification and claimed competitive performance on the GLUE benchmark. We reproduce their results on 3 GLUE datasets and extend them to 2 new datasets.

Methodology – We used the author’s code with some additions to accommodate our new datasets and other models. We ran all experiments using a combination of hyperparameters given by the original authors over 5 seeds, mainly using RTX8000 for 125 hours.

Results – We validate the original paper’s claims by reproducing its metrics within the reported standard deviation, proving the method’s competitiveness. Our extended trials highlight the method’s potential applicability to real-world data and reveal new considerations about prompt template, language model, and seed for optimal performance.

What was easy – The complete code implementation was available publicly with step-by-step instructions on how to get AMuLaP running with GLUE datasets. This made it easy to begin reproducing the work. Finding flagged and non-flagged tweets from Donald Trump to create our dataset was easy thanks to thetrumparchive.com which compiled this.

What was difficult – The original code is lengthy and lacks information on implementation dependencies, making it time-consuming to understand. It is made only for GLUE tasks and RoBERTa models and needs modification to work on new experiments. Some language models are not built for mask completion, thus not directly suitable for AMuLaP, and required us to search for solutions extensively. Finally, as label engineering is also tied to prompt engineering, we find expanding the work through the template given in the code challenging without prior experience in manual prompt engineering.

Communication with original authors – We did not feel the need to reach out to the authors as the method was well explained and simple to understand with basic probability knowledge and going through the code.

Copyright © 2023 V. Livernoche and V. Sujaya, released under a Creative Commons Attribution 4.0 International license.

Correspondence should be addressed to Victor Livernoche (victor.livernoche@mail.mcgill.ca)

The authors have declared that no competing interests exist.

Code is available at <https://github.com/vicliv/AMuLaP-Reproduction>. – SWH swh:1.dir:44eec5466ce396c0e4af4cba3e28b994c4a82600.

Open peer review is available at <https://openreview.net/forum?id=t8ZZ2Y356lx>.

1 Introduction

Prompt-based learning with pre-trained large language models (PLM) such as GPT-3 [2] have recently gained popularity and much work has been done in searching for prompting templates (i.e. prompt-engineering) [3] that best allow the PLM to treat tasks as masked language modeling problems. But when we want the PLM to fill in masks from a finite set of label words, which is a different approach to text classification, finding good label word sets is just as important (i.e. label engineering). The search for label words usually requires task expertise, knowledge of the internal PLM, or enough computation for a brute-forced search of the PLM’s vocabulary.

To mitigate the need for this manual work and expense, Wang, Xu, and McAuley^[1] proposed a novel technique for automatic multi-label map search through prompting. Their proposed solution, AMuLaP, can be described as using prompting on a set of same-class samples to find a group of label representatives for the class without needing human effort. We want to pick the best label to get the best classification for a specific language model. For example, “great” and “perfect” are two labels that can be associated with a positive sentiment and “awful” and “terrible” to a negative sentiment. Given these labels, AMuLaP finds the probability given by the model to complete the mask for each label and uses them for classification. The method is summarized in Figure 1 provided by the original paper. Using a group of labels reduces the noise brought by single labels when making predictions. On a high level, the authors claim competitive accuracy results on the GLUE dataset without fine-tuning nor access to model weights against majority voting and manual label zero-shot baselines. Its best performance is reached with fine-tuning, outperforming direct few-shot fine tuning on all GLUE datasets except for CoLA.

Automatic label search using prompt-based learning is relatively new, with only two works parallel to AMuLaP preceding it, called AutoL [4] and PETAL [5]. In this work, we aim to verify their main claims by reproducing their results on 3 GLUE datasets with and without fine-tuning and test their ability to generalize to real-world tasks and other models.

2 Scope of Reproducibility

2.1 Context and Original Claims

AMuLaP alleviates the need for manual label word search to optimize prompt-based learning with PLMs. Given a prompt template, AMuLaP uses a statistics-based selection method to choose multiple words to map as labels for a single class. As a result, the method’s way of finding label mappings brings the advantage of being parameter-free and simple. The authors claim AMuLaP’s ability to achieve competitive performance on GLUE tasks without accessing model weights, beating majority-vote and manual label zero-shot baselines across all GLUE datasets. Further, it reaches its best results with a fine-tuned model, outperforming direct few-shot fine-tuning on all GLUE datasets except for CoLA. AMuLaP’s works as its one-to-many label mapping per class reduce the noise brought by limited few shot examples and PLM complexities during automatic label selection and inference. Finally, they suggest that the increase of supervision pairs when selecting multiple labels during fine-tuning resembles working with augmented data, possibly being another reason to why multi-labels work.

2.2 Experiment Scope

Our work verifies the following claims of the paper and expands on it by exploring its generalization to real-world tasks and models. We use the original code, with modifications to accommodate two new non-GLUE datasets and models, and outline our experiments below:

- **GLUE performance reproduction:** We run AMuLaP under 2 settings (fine-tuned model and not) for MNLI, MNLI-mm, SST-2 and CoLA to verify claims of competitive performance on GLUE datasets.
- **Performance on 2 real-world datasets:** We run AMuLaP on 2 new datasets (flagged and non-flagged Trump tweets and Reddit comment-reply pairs). This is done for 3 goals: (1) to test the method’s robustness on unclean real-world data. (2) Verify the author’s claims related to multi-labels and noise by examining generated multi-labels. (3) explore the possible use case of AMuLaP’s multi-labels as a way to use PLMs to uncover insights about why datasets are classed a certain way. Particularly, we are interested in the multi labels as additional context behind why groups of Trump tweets were flagged, and as a clearer reason behind why a Reddit comment-reply pair falls in the “unsure” relation class.
- **Testing different language models:** We try AMuLaP with PLMs other than its RoBERTa-Large backbone and new manual templates to test the claim that the method works due to its use of multi-labels. In particular, we want to see whether the method’s success relies heavily on the choice of PLM/template rather than the idea of one-to-many label mappings.
- **Code modification:** The original implementation builds upon AutoL’s Gao, Fisch, and Chen^[4] codebase, a related work preceding it. To improve its usability, we modified the codebase to be more modular and added information on how to use it on non-GLUE datasets.

The rest of our report proceeds as follows. We first go over related work in section 3. This, section 4 outlines our methodology, and we summarize how AMuLaP works in section 4.1. Section 5 shows our results. Section 6 is a discussion of our insights from reproducing the work.

3 Related work

While in the past year, the principal method for learning natural language processing models was to pre-train them and fine-tune on a specific task, in the more recent years, there has been a shift with the use of prompting on pretrained models to predict masked tokens as explained by Liu et al. [3] The original GPT3 paper [2] had few-shot learning as the main focus; they showed that the model could use in-context learning from a few examples to complete a text. This further encouraged prompt engineering for classification tasks, i.e., converting a task into a language model format to use with a language model. This method was first explored by Trinh and Lee in 2018 [6]. They substituted a pronoun with one of two words with the most probable sentence for a language model. While simple, this new method started using prompt engineering for text classification since it achieved excellent results. Shin et al. [7] proposed a method for creating automatic prompting (AUTOPROMPT). They showed that we could use the masked language model “knowledge” without additional parameters or fine-tuning to perform sentiment analysis and natural language inference. This paper combines the original task inputs with a collection of trigger tokens learned using a gradient-based search method.

In the original AMuLaP paper [1], the authors argue that finding optimal label sets is

just as important as prompt template engineering to carry out prompting methods successfully on language models. Finding label maps would usually require task-specific expertise, and we can use the language model (LM) knowledge to find which words are “understood” well enough by the LM to be used as a label. Moreover, exhaustive trial and error approaches of all possible label mappings is intractable[4]. Automated Label Mapping approaches emerged as a solution to minimize the need for human involvement when selecting labels. The paper refers to two related works, and no other works have been released on automatic label mapping search between AMuLaP’s publication on April 2022 and December 2022 as far as we know.

Schick et al. [8] proposed PETAL (Pattern-Exploiting Training with Automatic Labels) in 2020. This framework builds on one of their earlier works, PET (Pattern-Exploiting Training), which has a verbalizer component that maps a label to a token representing its meaning. The Automatic Label aspect of PETAL searches for optimal verbalization candidates (label meaning representative tokens) using maximum-likelihood-based approaches. They also recognized the need for multiple verbalizers of a label for certain tasks (aligned with AMuLaP’s statement regarding the need for multi-labels). They implemented a solution that maps a single label to a set of tokens.

In 2021, Gao et al. [4] used a variation of Schick’s idea of automatic verbalizer search to find the optimal label words given a fixed prompting template. Their work varies from Schick’s by relying on a more brute-force method of specifying label word search space and selecting the words. In particular, within the top-k conditionally likely words within a language model’s vocabulary (where the condition is the input text), they search for the top n words that maximize zero shot accuracy of their training data and fine-tune on this n-sized set. Another variation comes in their final step of re-ranking the top n words using a development set. An important finding from both papers is that a good mapping from the original task labels to tokens is important to few-shot performance, which is why Wang et al. [1] focused on label engineering in their work.

4 Methodology

We use the author’s code with some modifications to handle non-BERT-based models and our two non-GLUE datasets. After these modifications, we follow the instructions in the README file to run experiments. We also use code from Gao, Fisch, and Chen^[4] to run an alternative automatic label search technique called AutoL on the Trump tweets and Debagreement. AutoL is an alternative to AMuLaP used by the authors as a comparison, which automatically finds single labels using prompts.

4.1 Model descriptions

AMuLaP is a method that uses the whole vocabulary of the language model to find the best labels. Using n examples per class, we sum the predictions of the mask token over the whole vocabulary. Then, for each word, we assign it to the class with the highest probability. Finally, we select k labels with the highest probability for each class. When testing, we sum the prediction on the mask token of these multi-labels for each class to get the prediction class with the highest sum. Refer to the original paper and Figure 1 for a formal description of the method [1].

The main model used is RoBERTa Large (355M parameters) [9] with pretrained weights modified to do prompt fine-tuning. In Table 4, we also used DerBERTa-V2 xlarge (900M) and xxlarge (1.5B) [10], RoBERTa Base (125M) [9], and BERT Large Uncased (340M) [11] for comparison. Similarly, the models were modified for prompt fine-tuning, and pretrained weights were used. From the original paper of each model, the results on the

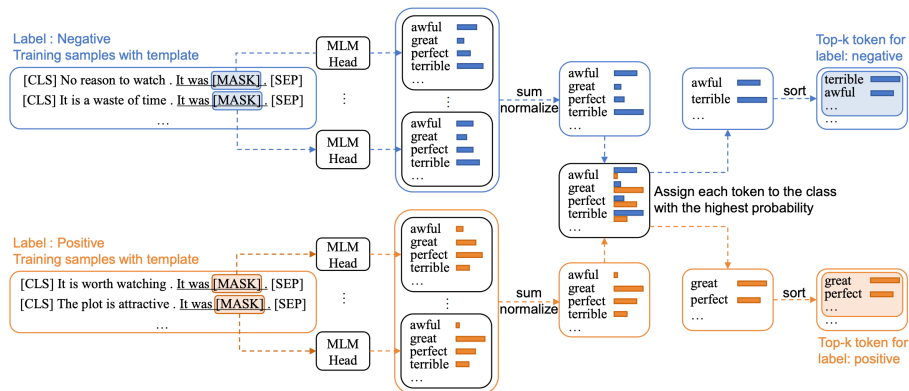


Figure 1. Illustration from the original paper of implementing AMuLaP on a binary sentiment classification task (SST-2). Each training sample with the task-specific template (the underlined text) is fed into a pre-trained language model L to get its own probability distribution over the vocabulary V . The obtained probability distributions are summed by class and normalized to get the probability distribution of each class. Then each token in V is assigned to the class with the highest probability (e.g., the token terrible is assigned to the class negative, and the token great to the class positive). Finally, we choose the top- k tokens as label words for each class [1]

GLUE benchmark [12] should be better the more parameters it has, apart from RoBERTa Base, which performs similar to and even outperforms BERT Large on certain tasks. A major limitation when choosing a model is that it has to have been trained with a mask token; otherwise we would need to fine-tune it, which might be resource expensive. To use a sequence-to-sequence model, we would need to modify the work extensively, and there is no say in whether the method would work as well.

4.2 Datasets

We use 4 GLUE datasets (MNLI, MNLI-mm, SST-2, and CoLA), a set of Donald Trump tweets from The Trump Tweet Archive (<https://www.thetrumparchive.com>) and Reddit comment reply pairs from Disagreement (<https://scale.com/open-av-datasets/oxford>)[13]. All of them were processed into train, dev (validation), and test splits using the author’s code to create k -shot data without any additional processing as we want to test the capacity of the method on “unclean” data. The code creates train and test files with n samples for each class, and we used $n = 16$ in our experiments. The test set is the dev set of an original GLUE dataset and the new dev set is the same length as the new train set. Below are details of each dataset:

- **MNLI:** MNLI is a dataset composed of sentence pairs annotated with textual entailment information (Entailment, contradiction, neutral) with each class represented approximately equally [14]. Here is an example:

Premise: Your gift is appreciated by each and every student who will benefit from your generosity.
Hypothesis: Hundreds of students will benefit from your generosity.
Label: neutral

- **SST-2:** The Stanford Sentiment Treebank regroups sentences labeled with a positive or negative sentiment. There are about the same amount of examples for each class [15]. Here is an example:

Positive example: a smile on your face

Negative example: the action is stilted

- **CoLA:** The Corpus of Linguistic Acceptability is a set of 10,657 English sentences labeled as grammatical or ungrammatical (acceptable or unacceptable) represented with 70.5% acceptable [16].

Example of an acceptable sentence: They made him president.

Example of an unacceptable sentence: The car honked down the road.

- **Trump Tweets:** Trump’s tweets are either flagged or not flagged by Twitter; these are our classes. Since there were only 304 flagged tweets out of 50,000 tweets, we created a dataset using randomly chosen 694 “ok” tweets to get 1000 data points with an average of 38 tokens per tweet. From now on, “ok” tweets refer to tweets that were not flagged by Twitter. We split the training with a ratio of 0.3, and the dev set is 0.1 of the test set.

Here is an example of a flagged tweet: RIGGED ELECTION. WE WILL WIN!

Here is an example of an ok tweet: So great to watch this! <https://t.co/pYoiLjM0pz>.

- **Disagreement:** It is a dataset composed 42,894 comment-reply pairs extracted from Reddit that are either agreeing, disagreeing, neutral or unsure [13]. We removed the unsure examples for our testing. Each class is represented equally in the dataset. On average, each pair consists of 89 tokens (using RoBERTa-Large tokenizer).

Example of a disagreement: Comment: I feel like the pictures should be reversed... Only one side is throwing a fit.

Reply: We’ve had months of leftists burning and looting, you think the right is throwing a fit?

4.3 Hyperparameters

We replicated the use of hyperparameters in the original paper as much as possible; we only tried one learning rate and one batch size instead of 3 as running the method when fine-tuning was computationally expensive. For all GLUE experiments, we tried five different seeds (13, 21, 42, 87, and 100) to generate the training split. For Trump and Disagreement, we only used a single set of training split. We used the same seeds for all experiments to get the standard deviation and average. We report results from the best-performing top-k parameter based on the accuracy of the dev set. We used top-k’s of 1, 2, 4, 8, and 16 for fine-tuning, and 1, 2, 4, 8, 16, 32, 64, 128, 256, 512, and 1024 for no fine-tuning. We used a batch size of 8 for all experiments and a learning rate of 1e-5 when fine-tuning. Figure 2 shows the impact of the top-k on SST-2 and MNLI-mm. We found that depending on the template and the dataset, the best top-k can be large, but most of the time, there are no specific patterns to determine the best top-k. Larger top-k reduces the standard error on the results and with a good template, it should perform well compared to a small one.

4.4 Experimental setup and code

The experiment can be run using a python file with arguments. We ran each dataset with the hyperparameters described in Section 4.3 to report the best result based on the dev set. A file named `prompt.py` has a class containing multi-label prompting methods. New models that support mask modeling can easily be added in `model.py` by replicating the code used for the other classes. Every dataset was evaluated using accuracy, i.e. the percentage of correctly classified examples, but CoLA uses Matthews Correlation which

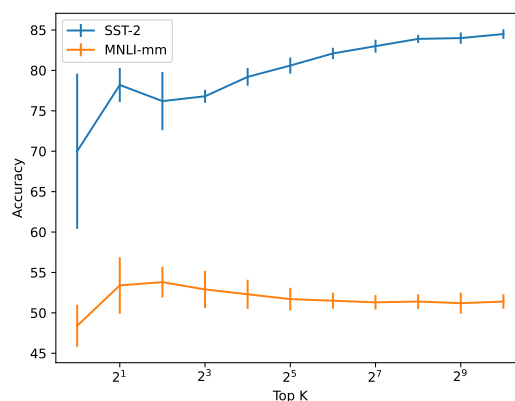


Figure 2. Accuracy on SST-2 and MNLI-mm using AMuLaP with different top K's with no fine-tuning.

is a metric that is used for imbalanced classes that is related to the F1-score (1 for perfect agreement, -1 for perfect disagreement). The code is available at <https://github.com/vicliv/AMuLaP-Reproduction>. Refer to the README.md file for more information to run the code.

4.5 Computational requirements

We primarily used RTX8000 for AMuLaP runs, which took 2 minutes per run without fine-tuning, and 10 minutes per run with fine-tuning on average. We took an average of 250 minutes for fine-tuning experiments and 110 minutes non fine-tuned experiments. When working with DeBERTa-V2-xxlarge, we used A100 with 80 GB of memory. AutoL experiments were run on standard Google Colab and Kaggle GPU, with single runs taking 2 and 10 hours for the Trump and Deagreement datasets, respectively. In total, with the ability to run over some experiments in parallel, our main AMuLaP experiments took a total of 125 hours.

5 Results

Our results generally support the claim made by the original authors, of AMuLaP's competitive results on the GLUE datasets and the reduction of noise by using multiple labels. Further, our experiments on Trump tweets and Deagreement showed AMuLaP's potential in generalizing to real-world data. However, we found a large decline in performance when using non-default PLMs and templates.

5.1 Results reproducing original paper

We ran the experiments on SST-2, MNLI, MNLI-mm, and CoLA to reproduce the results from the original paper (first point in section 2.2), the two of which are compared in Table 1. The average results we got for AMuLaP without fine-tuning is within the standard deviation of the original work and slightly worse when fine-tuning. Overall, the results correspond to what was shown in the original paper.

We got slightly worse results when fine-tuning, which might be because we only tried a single learning rate. Trying multiple learning rates would have been too long and would have use much more computation resources.

	MNLI (acc)	MNLI-mm (acc)	SST-2 (acc)	CoLA (Matt.)
Setting 2: Dtrain + Ddev ; No parameter update.				
Original AMuLaP	50.8 ± 2.1	52.3 ± 1.8	86.9 ± 1.6	2.3 ± 1.4
Reproduced AMuLaP	51.3 ± 1.9	52.9 ± 2.3	86.9 ± 1.4	2.2 ± 2.9
Setting 3: Dtrain + Ddev ; Prompt-based fine-tuning.				
Original AMuLaP FT	70.6 ± 2.7	72.5 ± 2.4	93.2 ± 0.7	18.3 ± 9.4
Reproduced AMuLaP FT	67.0 ± 2.7	69.1 ± 3.2	92.6 ± 1.4	17.0 ± 11.1

Table 1. Reproduced and original experiment results under two settings on some GLUE dataset using Roberta. For few-shot settings, n is set to 16 per class. For reproduction, we give the average and standard deviation over 5 runs (different seeds) for all experiments.

5.2 Results beyond original paper

Default AMuLaP on Trump tweets and Debagreement – Table 2 shows the results of running AMuLaP with RoBERTa-Large for the Trump and Debagreement datasets. We also ran AutoL to use as a comparison, similar to the original paper. We find relatively high accuracies for Trump: 77.3% without fine-tuning and 86.7% with, but unsatisfactory accuracies for Debagreement: 30.5% without fine-tuning and 40.5% with.

	Trump (acc)	Debagreement (acc)
Setting 2: Dtrain + Ddev ; No parameter update.		
AMuLaP	77.3	30.5
Setting 3: Dtrain + Ddev ; Prompt-based fine-tuning.		
AMuLaP FT	86.7 ± 2.7	40.5 ± 1.3
Auto-L FT	89.1 ± 1.1	36.4

Table 2. Experiment results under two settings for the Trump and Debagreement datasets using AMuLaP (multi-label) and AutoL (single label) with Roberta Large. For few-shot settings, k is set to 16 per class. We show the average and standard deviation over 5 runs with different seeds for all experiments apart from Auto-L on Debagreement, done on a single run.

AMuLaP for Trump tweets and Debagreement: Label and Template – Table 3 shows the multi-labels found using AMuLaP for the Trump and Debagreement datasets along with the fixed manual templates used to find them:

The table shows a higher amount of noise (non-word labels) chosen as multi labels for the Debagreement dataset compared to Trump tweets, which in fact, presents none.

AMuLaP with Alternative PLMs – Following point 3 in section 2.2, we test AMuLaP with PLMs other than its default RoBERTa-Large backbone and report results in Table 4. In this experiment, we take the results using the backbone RoBERTa-Large as a baseline. We observe that the baseline model performs best for the Debagreement dataset in both settings and SST-2 in the fine-tuned setting. It is outperformed in the non-fine-tuned setting only by DeBERTa for the SST-2 and Trump datasets and in the fine-tuned setting by RoBERTa-Base for the Trump dataset.

6 Discussion

Our work verifies the authors’ claims about AMuLaP on the GLUE dataset, as our reproduced results fall within the standard deviation of the original report. Testing AutoL

	Generated Multi-Labels
Trump Template 1: "Sentence 1. It was [MASK]."	
Flagged tweets	rigged, awesome, close, disgusting, terrible, huge, over, shocking, illegal, ridiculous
Ok tweets	fun, true, amazing, inevitable, not, great, beautiful, funny, good, incredible
Trump Template 2: "Sentence 1. This is a [MASK] tweet."	
Flagged tweets	false, TRUE, warning, lawful, FALSE, fraud, rigged, harassing, felony, panicked
Ok tweets	recent, real, new, live, deleted, related, verified, true, great, sponsored
Debate Template 1: "Sentence 1. [MASK]. Sentence 2."	
Disagreeing pairs	Or, And, </s>, Mattis, Something, Like, *, Even, Wait, -
Agreeing pairs	But, The, If, and, Now, New, Assuming, US, Yes, Heck
Neutral pairs	So, It, Well, Maybe, Just, ., Also, <, Obviously, ...
Unsure pairs	Oh, Why, I, Because, As, but, Ok, However, Not, OK

Table 3. Multi labels found for each dataset class using a specific fixed manual template with AMuLaP and RoBERTa-Large

	SST-2 (acc)	Trump (acc)	Debate (acc)
Setting 2: Dtrain + Ddev ; No parameter update.			
AMuLaP RoBERTa-Large (Baseline)	86.9 ± 1.6	77.3	30.5
AMuLaP Deberta-V2-xxlarge	89.3 ± 0.7	78.9	27.5
AMuLaP Deberta-V2-xlarge	81.9 ± 2.4	79.4	28.7
AMuLaP Roberta-Base	84.5 ± 0.6	63.0	25.0
AMuLaP Bert-Large-Uncased	73.8 ± 2.0	63.0	26.6
Setting 3: Dtrain + Ddev ; Prompt-based fine-tuning.			
AMuLaP RoBERTa-Large (Baseline)	93.2 ± 0.7	86.7 ± 2.7	40.5 ± 1.3
AMuLaP Roberta-Base	89.5 ± 0.4	76.7 ± 1.7	39.0 ± 0.6
AMuLaP Bert-Large-Uncased	84.7 ± 1.6	82.9 ± 2.4	30.6 ± 0.5
Auto-L Roberta-Base	-	80.3 ± 2.1	31.1

Table 4. Experimental results under two settings. Reproduction of results from using different models on SST-2, Trump, and Debate datasets for AMuLaP. For few-shot settings, n is set to 16 per class. We show the average and standard deviation over 5 runs for all experiments apart from Auto-L on Debate, which was done on a single run.

on Trump and Debate also supports the author’s motivation for using multiple labels. Instead of multi-labels, Auto-L searches for good combinations of single labels per class and consistently achieves results lower than AMuLaP with the same PLM.

Our extended experiments uncover the potential use of AMuLaP and the factors needed for its success. AMuLaP’s multi-labels can be used to find possible reasons behind classification, as common interpretations of the Trump multi-labels make sense as the reasoning behind why tweets were flagged (labeled “false” or “rigged”) or not (labeled “real” or “true”). Then, examining our poor Debate results uncovered a need to explore effective templates to support automatic label search. Table 4 shows PLM choices impact performance. Also, the 0.48 average standard deviation increase in our reproduction of GLUE dataset results in table 1 shows the significance of seeds in AMuLaP’s performance. We elaborate on our observations regarding template significance and PLM choice below.

In the classification task lens, AMuLaP’s high accuracy on the Trump dataset is proof of its potential on real-world data, while its poor results on Debagreement raise new considerations. A closer look at the generated labels shows AMuLaP’s fixed template aspect as a weakness. The generated labels (shown in table 3) contain noise, and 90% of them are conjunction words. The selection of these words makes sense as the template places the mask token between the comment and reply pair. But, the occurrences of similar purpose conjunctions (ex: “Even”, “But”, “Heck” and “However” across three different classes) suggest that they may not be informative enough for accurate classification.

Informed by this observation, we tested other manual templates that resemble successful Trump ones by putting the mask token at the prompt end. However, we could not find a template that reaches a higher accuracy. We attempted using the automatic template search functionalities of AutoL [4], but needed more resources to complete it. Joint automatic template and label search remain a future task for us.

Using different PLMs, table 4 verify the authors’ claim that AMuLaP achieves its best results with access to the model weights. All experiments under setting 3 outperform their equivalent in setting 2. We observe that larger PLMs do not consistently achieve better results, as shown by the baseline RoBERTa-Large model exceeding DeBERTa-V2-xlarge for SST-2 in setting 2. We are interested in how well AMuLaP can use larger models; as we have seen, DeBERTa-V2 achieved better results than RoBERTa while still being relatively small compared to OPT models of GPT-3. We wanted to use GPT-3, but it is currently impossible to get the predictions over the whole vocabulary as AMuLaP requires. This feature may be available in the future, or we will need to contact OpenAI.

For future work, finding more appropriate baselines to measure our findings on Trump tweets and Debagreement would be helpful. Using human evaluation, where literature professors and student judge whether the automatically searched multiple labels are good representations of a class, would be a good baseline to improve upon.

6.1 What was easy

The complete code implementation was available publicly with step-by-step instructions on how to get AMuLaP running with GLUE datasets. This made it easy to begin the reproduction of the work. Finding flagged and not flagged tweets from Donald Trump was also easier than we initially thought. We were able to easily create our dataset thanks to thetrumparchive.com which compiled the tweets and labeled each of them.

6.2 What was difficult

We needed to understand the code before fully understanding the explanation of AMuLaP in the original paper. However, the code was built upon another paper’s work and contained lengthy files, making it time-consuming to understand. Further, the implementation works only with GLUE datasets and BERT-based PLMs, without adequate documentation on what to change to work with new data and non-BERT PLMs. Most models were not trained with a mask token and, thus, were unsuitable for a prompt-based method with a mask, such as how AMuLaP was implemented. As we wanted to try a larger model than RoBERTa, we had to search to find a way to find one extensively. We used DeBERTa-V2 for mask modeling, but the most recent version of Hugging-Face transformers had an issue loading the weight of the model head. Fortunately, the issue was being solved concurrently with our research, and we could use a branch directory. Because of time constraints and GPU memory limitations, we could not fine-tune DeBERTa-V2. Finally, the label engineering problem AMuLaP tries to solve is tied to prompt engineering, so we find template experimentation challenging without prior experience.

References

1. H. Wang, C. Xu, and J. McAuley. "Automatic Multi-Label Prompting: Simple and Interpretable Few-Shot Classification." In: **Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies**. Seattle, United States: Association for Computational Linguistics, July 2022, pp. 5483–5492. doi: 10.18653/v1/2022.naacl-main.401. URL: <https://aclanthology.org/2022.naacl-main.401>.
2. T. B. Brown et al. **Language Models are Few-Shot Learners**. 2020. doi: 10.48550/ARXIV.2005.14165. URL: <https://arxiv.org/abs/2005.14165>.
3. P. Liu, W. Yuan, J. Fu, Z. Jiang, H. Hayashi, and G. Neubig. **Pre-train, Prompt, and Predict: A Systematic Survey of Prompting Methods in Natural Language Processing**. 2021. doi: 10.48550/ARXIV.2107.13586. URL: <https://arxiv.org/abs/2107.13586>.
4. T. Gao, A. Fisch, and D. Chen. "Making Pre-trained Language Models Better Few-shot Learners." In: **Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)**. Online: Association for Computational Linguistics, Aug. 2021, pp. 3816–3830. doi: 10.18653/v1/2021.acl-long.295. URL: <https://aclanthology.org/2021.acl-long.295>.
5. T. Schick, H. Schmid, and H. Schütze. "Automatically Identifying Words That Can Serve as Labels for Few-Shot Text Classification." In: **Proceedings of the 28th International Conference on Computational Linguistics**. Barcelona, Spain (Online): International Committee on Computational Linguistics, Dec. 2020, pp. 5569–5578. doi: 10.18653/v1/2020.coling-main.488. URL: <https://aclanthology.org/2020.coling-main.488>.
6. T. H. Trinh and Q. V. Le. **A Simple Method for Commonsense Reasoning**. 2018. doi: 10.48550/ARXIV.1806.02847. URL: <https://arxiv.org/abs/1806.02847>.
7. T. Shin, Y. Razeghi, R. L. Logan, E. Wallace, and S. Singh. **AutoPrompt: Eliciting Knowledge from Language Models with Automatically Generated Prompts**. 2020. doi: 10.48550/ARXIV.2010.15980. URL: <https://arxiv.org/abs/2010.15980>.
8. T. Schick, H. Schmid, and H. Schütze. "Automatically identifying words that can serve as labels for few-shot text classification." In: **arXiv preprint arXiv:2010.13641** (2020).
9. Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis, L. Zettlemoyer, and V. Stoyanov. **RoBERTa: A Robustly Optimized BERT Pretraining Approach**. 2019. doi: 10.48550/ARXIV.1907.11692. URL: <https://arxiv.org/abs/1907.11692>.
10. P. He, X. Liu, J. Gao, and W. Chen. "DEBERTA: DECODING-ENHANCED BERT WITH DISENTANGLED ATTENTION." In: **International Conference on Learning Representations**. 2021. URL: <https://openreview.net/forum?id=XPZlaotutsD>.
11. J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova. "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding." In: **Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)**. Minneapolis, Minnesota: Association for Computational Linguistics, June 2019, pp. 4171–4186. doi: 10.18653/v1/N19-1423. URL: <https://aclanthology.org/N19-1423>.
12. A. Wang, A. Singh, J. Michael, F. Hill, O. Levy, and S. Bowman. "GLUE: A Multi-Task Benchmark and Analysis Platform for Natural Language Understanding." In: **Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP**. Brussels, Belgium: Association for Computational Linguistics, Nov. 2018, pp. 353–355. doi: 10.18653/v1/W18-5446. URL: <https://aclanthology.org/W18-5446>.
13. J. Poug  -Biyong, V. Semanova, A. Matton, R. Han, A. Kim, R. Lambiotte, and D. Farmer. "DEBAGREEMENT: A comment-reply dataset for (dis) agreement detection in online debates." In: **Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)**. 2021.
14. A. Williams, N. Nangia, and S. Bowman. "A Broad-Coverage Challenge Corpus for Sentence Understanding through Inference." In: **Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)**. New Orleans, Louisiana: Association for Computational Linguistics, 2018, pp. 1112–1122. URL: <http://aclweb.org/anthology/N18-1101>.
15. R. Socher, A. Perelygin, J. Wu, J. Chuang, C. D. Manning, A. Ng, and C. Potts. "Recursive Deep Models for Semantic Compositionality Over a Sentiment Treebank." In: **Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing**. Seattle, Washington, USA: Association for Computational Linguistics, Oct. 2013, pp. 1631–1642. URL: <https://www.aclweb.org/anthology/D13-1170>.
16. A. Warstadt, A. Singh, and S. R. Bowman. "Neural Network Acceptability Judgments." In: **Transactions of the Association for Computational Linguistics 7** (2019), pp. 625–641. doi: 10.1162/tacl_a_00290. URL: <https://aclanthology.org/Q19-1040>.

A Dataset examples

We give a few more examples from our social media datasets.

A.1 Trump's Tweets

- **Sentence:** ...What are they trying to hide. They know, and so does everyone else. EXPOSE THE CRIME! . (Flagged)
- **Sentence:** Immigration reform really changes the voting scales for the Republicans—for the worse! . (Not flagged)
- **Sentence :** The United States better address China's exchange rate before they steal our country and it is too late! China is laughing at us. . (Not flagged)
- **Sentence:** Corrupt Election! <https://t.co/MxNjfCEtKP> . (Flagged)

A.2 Debagreement

- **Comment:** Unfortunately I think they'll just move the voting to be virtual and proceed with the timeline they had planned.
Reply: So they cant. McConnell blocked ... (Disagree)
- **Comment:** This will annoy the hell out of the liberals who regularly use Uber and Lyft. I see this as a good thing!
Reply: Looks like they pushed all their chips in the pot with a pair of deuces and it wasn't enough. LOL (Neutral)
- **Comment:** Perhaps the older generation want a return to the blitz, the time has warped their recollection of rationing.
Reply: I asked grandpa about the Blitz once! Indeed he said that rationing was important. (Agree)
- **Comment:** Why is that bad news? I hope other addicts will see that and know there is hope to get out.
Reply: Am I being down voted for suggesting that its a good thing for drug addicts recover? (Unsure)

B Labels Comparison between Models

We show the labels chosen using AMuLap for Roberta-large and BERT-base for the SST-2 dataset. It is interesting to see the difference between each model. It shows that they put importance on different words.

RoBERTa

- **positive:** Great, perfect, fun, brilliant, amazing, good, wonderful, beautiful, excellent, fantastic
- **Negative:** Terrible, awful, disappointing, not, horrible, obvious, funny, inevitable, bad, boring

BERT

- **Positive:** perfect, fun, beautiful, good, great, amazing, wonderful, incredible, fantastic, successful, magnificent, awesome, nice, easy, spectacular, glorious
- **Negative:** horrible, true, awful, funny, stupid, brilliant, wrong, terrible, me, ridiculous, right, simple, crazy, there, real, him