

# Comparative Analysis of Prompt Engineering Techniques in Fine-Tuning Chatbots with skt/kogpt2-base-v2

Anonymous NLPAICS submission

## Abstract

Since OpenAI released GPT, the field of NLP and LLMs has experienced rapid advancements over the years, a trend that continues even today. Among these, prompt engineering has emerged as a key technique for effectively leveraging LLMs, becoming a critical area of global interest in language model development. Numerous papers and practical experiences have demonstrated the substantial impact and significance of prompts, leading companies and research institutions to make remarkable progress in this field. However, as long as the internal workings of the models remain partially opaque, there will always be room for improvement and exploration in this domain.

## 1 Introduction

This study focuses on a detailed exploration of the application of prompts. While the effects of prompt engineering in large-scale models have been validated across various dimensions, do these findings universally apply to all LLMs? Are they truly optimized and generalized effectively? Could there be a gap between general users' perspectives and the outcomes of prompt engineering? What differences exist between prompts used in fine-tuning processes and those utilized during the inference phase (e.g., by users or questioners)? Based on these questions, this research conducts experiments to address these issues.

## 2 Method

The experiments in this study focus on prompt engineering, with the assumption that its comparative utility would be more pronounced in LLMs exhibiting clear emergent abilities (models with sufficient parameters to demonstrate such capabilities). However, due to practical constraints, a compromise was made by selecting a more accessible model for experimentation. Specifically, the skt/kogpt2-

base-v2 model, pretrained on Korean corpora, was adopted for this research (see [GitHub repository](#)).

### 2.1 2-1 Data

Data was collected from four datasets, provided on the KoChatGPT GitHub repository, totaling 167,577 question-answer pairs:

- data1: ChatbotData (11,824 sentences)
- data2: AI Hub Korean Conversations (49,711 sentences)
- data3: AI Hub General Knowledge (100,268 sentences)
- data4: KorQuAD (5,774 sentences)

From this, a dataset of 12,000 randomly sampled questions was primarily used.

Additionally, the entire KorQuAD 2.1 dataset was utilized, which consists of 47,957 Wikipedia articles and 102,960 question-answer pairs. These pairs are divided into a training set (83,486 pairs) and a development (dev) set (10,165 pairs). A subset of the training dataset, consisting of 48,111 pairs, was used for this study. The dataset includes an official evaluation script, input sample predictions, and a leaderboard for model evaluation.

However, as this study requires internal comparative analysis, the dev dataset was repurposed as the test set for evaluation. The dataset and resources are available on the [KorQuAD GitHub repository](#).

### 2.2 2-2 Model

The model used in this study, skt/kogpt2-base-v2, is based on the GPT-2 architecture, a transformer-based model designed for natural language processing tasks. The architecture utilizes a multi-head self-attention mechanism and a feed-forward neural network, consisting of 12 layers, 12 attention heads, and a hidden size of 768. This design allows the model to capture complex dependencies within text, making it highly effective for tasks involving contextual understanding.

skt/kogpt2-base-v2 is specifically pretrained on

Korean corpora, enabling it to understand and generate text in Korean effectively. The pretraining involved extensive datasets, including conversational, general knowledge, and domain-specific texts, to enhance its adaptability to diverse Korean linguistic contexts. This pretraining allows the model to excel in Korean-specific syntax, grammar, and semantic nuances.

This model’s ability to be fine-tuned for downstream tasks, such as question answering or text generation, ensures its versatility. By leveraging this architecture, the study focuses on optimizing prompt engineering techniques tailored for Korean language data.

### 2.3 2-3 Train

A total of 47,970 question-answer pair data points will be used for the chatbot downstream task, and the model’s results and comparisons will be evaluated using a test dataset of 10,165 question-answer pairs. The answers in the test dataset will be compared with the responses generated by the trained model based on the test questions.

| Command             | Output |
|---------------------|--------|
| <code>{\`a}</code>  | ä      |
| <code>{\^e}</code>  | ê      |
| <code>{\`i}</code>  | ì      |
| <code>{\ .I}</code> | İ      |
| <code>{\o}</code>   | ø      |
| <code>{\`u}</code>  | ú      |
| <code>{\aa}</code>  | å      |

Example commands for accented characters, to be used in, *e.g.*, BibTeX entries.

## 3 Preamble

The first line of the file must be

```
\documentclass[11pt]{article}
```

To load the style file in the review version:

```
\usepackage[review]{NLPAICS2024}
```

For the final version, omit the review option:

```
\usepackage{NLPAICS2024}
```

To use Times Roman, put the following in the preamble:

```
\usepackage{times}
```

(Alternatives like txfonts or newtx are also acceptable.) Please see the L<sup>A</sup>T<sub>E</sub>X source of this

document for comments on other packages that may be useful. Set the title and author using `\title` and `\author`. Within the author list, format multiple authors using `\and` and `\And` and `\AND`; please see the L<sup>A</sup>T<sub>E</sub>X source for examples. By default, the box containing the title and author names is set to a minimum of 5 cm. If you need more space, include the following in the preamble:

```
\setlength\titlebox{<dim>}
```

where `<dim>` is replaced with a length. Do not set this length smaller than 5 cm.

## 4 Document Body

### 4.1 Footnotes

Footnotes are inserted with the `\footnote` command.<sup>1</sup>

### 4.2 Tables and figures

See Table ?? for an example of a table and its caption. **Do not override the default caption sizes.**

### 4.3 Hyperlinks

Users of older versions of L<sup>A</sup>T<sub>E</sub>X may encounter the following error during compilation:

```
\pdfendlink ended up in different
nesting level than \pdfstartlink.
```

This happens when pdfL<sup>A</sup>T<sub>E</sub>X is used and a citation splits across a page boundary. The best way to fix this is to upgrade L<sup>A</sup>T<sub>E</sub>X to 2018-12-01 or later.

### 4.4 Citations

Table 1 shows the syntax supported by the style files. We encourage you to use the natbib styles. You can use the command `\citet` (cite in text) to get “author (year)” citations, like this citation to a paper by Gusfield (1997). You can use the command `\citep` (cite in parentheses) to get “(author, year)” citations (Gusfield, 1997). You can use the command `\citealp` (alternative cite without parentheses) to get “author, year” citations, which is useful for using citations within parentheses (e.g. Gusfield, 1997).

<sup>1</sup>This is a footnote.

| Output                        | natbib command   | Old ACL-style command |
|-------------------------------|------------------|-----------------------|
| (Cooley and Tukey, 1965)      | \citep           | \cite                 |
| Cooley and Tukey, 1965        | \citealp         | no equivalent         |
| Cooley and Tukey (1965)       | \citet           | \newcite              |
| (1965)                        | \citeyearpar     | \shortcite            |
| Cooley and Tukey’s (1965)     | \citeposs        | no equivalent         |
| (FFT; Cooley and Tukey, 1965) | \citep[FFT;]{}[] | no equivalent         |

Table 1: Citation commands supported by the style file. The style is based on the natbib package and supports all natbib citation commands. It also supports commands defined in previous ACL style files for compatibility.

## 4.5 References

The L<sup>A</sup>T<sub>E</sub>X and BibT<sub>E</sub>X style files provided roughly follow the American Psychological Association format. If your bib file is named custom.bib, then placing the following before any appendices in your L<sup>A</sup>T<sub>E</sub>X file will generate the references section for you:

```
\bibliographystyle{acl_natbib}
\bibliography{custom}
```

You can obtain the complete ACL Anthology as a BibT<sub>E</sub>X file from <https://aclweb.org/anthology/anthology.bib.gz>. To include both the Anthology and your own .bib file, use the following instead of the above.

```
\bibliographystyle{acl_natbib}
\bibliography{anthology,custom}
```

Please see Section 5 for information on preparing BibT<sub>E</sub>X files.

## 4.6 Appendices

Use \appendix before any appendix section to switch the section numbering over to letters. See Appendix A for an example.

## 5 BibT<sub>E</sub>X Files

Unicode cannot be used in BibT<sub>E</sub>X entries, and some ways of typing special characters can disrupt BibT<sub>E</sub>X’s alphabetization. The recommended way of typing special characters is shown in Table ?? . Please ensure that BibT<sub>E</sub>X records contain DOIs or URLs when possible, and for all the ACL materials that you reference. Use the doi field for DOIs and the url field for URLs. If a BibT<sub>E</sub>X entry has a URL or DOI field, the paper title in the references section will appear as a hyperlink to the paper, using the hyperref L<sup>A</sup>T<sub>E</sub>X package.

## Limitations

NLPAICS2024 requires all submissions to have a section titled “Limitations”, for discussing the limitations of the paper as a complement to the discussion of strengths in the main text. This section should occur after the conclusion, but before the references. It will not count towards the page limit. The discussion of limitations is mandatory. Papers without a limitation section will be desk-rejected without review.

While we are open to different types of limitations, just mentioning that a set of results have been shown for English only probably does not reflect what we expect. Mentioning that the method works mostly for languages with limited morphology, like English, is a much better alternative. In addition, limitations such as low scalability to long text, the requirement of large GPU resources, or other things that inspire crucial further investigation are welcome.

## Ethics Statement

Scientific work published at NLPAICS2024 must comply with the ACL Ethics Policy.<sup>2</sup> We encourage all authors to include an explicit ethics statement on the broader impact of the work, or other ethical considerations after the conclusion but before the references. The ethics statement will not count toward the page limit (8 pages for long, 4 pages for short papers).

## Acknowledgements

This document has been adapted by Jordan Boyd-Graber, Naoaki Okazaki, Anna Rogers from the style files used for earlier ACL, EMNLP and NAACL proceedings, including those for NLPAICS 2024 by Ignatius Ezeani and Ruslan

<sup>2</sup><https://www.aclweb.org/portal/content/acl-code-ethics>

|     |   |   |     |
|-----|---|---|-----|
| 223 | Mitkov, EACL 2023 by Isabelle Augenstein and              | Dan Gusfield. 1997. <i>Algorithms on Strings, Trees and</i>     | 273 |
| 224 | Andreas Vlachos, EMNLP 2022 by Yue Zhang,                 | <i>Sequences</i> . Cambridge University Press, Cambridge,       | 274 |
| 225 | Ryan Cotterell and Lea Frermann, ACL 2020 by              | UK.   | 275 |
| 226 | Steven Bethard, Ryan Cotterell and Rui Yan, ACL           | Mary Harper. 2014. <i>Learning from 26 languages: Pro-</i>      | 276 |
| 227 | 2019 by Douwe Kiela and Ivan Vulić, NAACL                 | <i>gram management and science in the babel program</i> . In    | 277 |
| 228 | 2019 by Stephanie Lukin and Alla Roskovskaya,             | <i>Proceedings of COLING 2014, the 25th International</i>       | 278 |
| 229 | ACL 2018 by Shay Cohen, Kevin Gimpel, and                 | <i>Conference on Computational Linguistics: Technical</i>       | 279 |
| 230 | Wei Lu, NAACL 2018 by Margaret Mitchell and               | <i>Papers</i> , page 1, Dublin, Ireland. Dublin City University | 280 |
| 231 | Stephanie Lukin, BibT <sub>E</sub> X suggestions for      | and Association for Computational Linguistics.                  | 281 |
| 232 | (NA)ACL 2017/2018 from Jason Eisner, ACL                  | Mohammad Sadegh Rasooli and Joel R. Tetreault. 2015.            | 282 |
| 233 | 2017 by Dan Gildea and Min-Yen Kan, NAACL                 | <i>Yara parser: A fast and accurate dependency parser</i> .     | 283 |
| 234 | 2017 by Margaret Mitchell, ACL 2012 by Maggie             | <i>Computing Research Repository</i> , arXiv:1503.06733.        | 284 |
| 235 | Li and Michael White, ACL 2010 by Jing-Shin               | Version 2.  | 285 |
| 236 | Chang and Philipp Koehn, ACL 2008 by Johanna              | <b>A Example Appendix</b>                                       | 286 |
| 237 | D. Moore, Simone Teufel, James Allan, and                 | This is a section in the appendix.                              | 287 |
| 238 | Sadaoki Furui, ACL 2005 by Hwee Tou Ng and                |   |     |
| 239 | Kemal Oflazer, ACL 2002 by Eugene Charniak                |   |     |
| 240 | and Dekang Lin, and earlier ACL and EACL                  |   |     |
| 241 | formats written by several people, including John         |   |     |
| 242 | Chen, Henry S. Thompson and Donald Walker.                |   |     |
| 243 | Additional elements were taken from the                   |   |     |
| 244 | formatting instructions of the <i>International Joint</i> |   |     |
| 245 | <i>Conference on Artificial Intelligence</i> and the      |   |     |
| 246 | <i>Conference on Computer Vision and Pattern</i>          |   |     |
| 247 | <i>Recognition</i> .                                      |   |     |

## References

Rie Kubota Ando and Tong Zhang. 2005. *A framework for learning predictive structures from multiple tasks and unlabeled data*. *Journal of Machine Learning Research*, 6:1817–1853.

Galen Andrew and Jianfeng Gao. 2007. *Scalable training of  $L_1$ -regularized log-linear models*. In *Proceedings of the 24th International Conference on Machine Learning*, pages 33–40.

Isabelle Augenstein, Tim Rocktäschel, Andreas Vlachos, and Kalina Bontcheva. 2016. *Stance detection with bidirectional conditional encoding*. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 876–885, Austin, Texas. Association for Computational Linguistics.

James W. Cooley and John W. Tukey. 1965. *An algorithm for the machine calculation of complex Fourier series*. *Mathematics of Computation*, 19(90):297–301.

James Goodman, Andreas Vlachos, and Jason Naradowsky. 2016. *Noise reduction and targeted exploration in imitation learning for Abstract Meaning Representation parsing*. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1–11, Berlin, Germany. Association for Computational Linguistics.