

1 LOST IN TRANSLATION: TRANSLATING GENERATION
2 ALPHA INTERNET SLANG USING MACHINE LEARNING

3 A Special Problem Proposal
4 Presented to
5 the Faculty of the Division of Physical Sciences and Mathematics
6 College of Arts and Sciences
7 University of the Philippines Visayas
8 Miag-ao, Iloilo

9 In Partial Fulfillment
10 of the Requirements for the Degree of
11 Bachelor of Science in Computer Science by

12 FLAUTA, Neil Bryan
13 GIMENO, Ashley Joy
14 GIMENO, Carl Jorenz

15 Francis DIMZON
16 Adviser

17 December 9, 2024

Abstract

19 Internet slang is an informal variation of language that is prominent to the younger
20 generation. The usage of this language brought generational divide between them
21 and the older generations. This study aims to develop a translation tool leveraging
22 Large Language Models (LLMs) to bridge this issue. A dataset of Generation
23 Alpha slang sentences and their formal equivalents will be used to fine-tune an
24 existing LLM. The model will be trained to translate slang sentences into formal
25 English, and will be compared against the baseline model using various evaluation
26 metrics. The study highlights the significance of addressing communication gaps
27 and provides insights into how technology can enhance understanding and reduce
28 miscommunications across generations. This research contributes to the broader
29 discourse on language adaptation and generational communication in the digital
30 age.

31 **Keywords:** Internet Slang, Generation Alpha, Generational Divide,
LoRA, LLM

Contents

33	1 Introduction	1
34	1.1 Overview	1
35	1.2 Problem Statement	2
36	1.3 Research Objectives	3
37	1.3.1 General Objectives	3
38	1.3.2 Specific Objectives	3
39	1.4 Scope and Limitations of the Research	3
40	1.5 Significance of the Research	3
41	2 Review of Related Literature	4
42	2.1 Communication Gap between Generations	4
43	2.2 Existing Studies	5
44	2.3 LoRA for Fine Tuning	6
45	2.4 Chapter Summary	6
46	3 Research Methodology	7
47	3.1 Research Activities	7
48	3.1.1 Creation of the dataset	7

49	3.1.2	Identification of potential LLM to be used	7
50	3.1.3	Lookup on available GPU on demand services	8
51	3.1.4	Study on LoRA implementation for LLM	8
52	3.1.5	Preprocessing of data	8
53	3.1.6	Prototype implementation of LoRA	9
54	3.1.7	Implementation of LoRA on selected model	9
55	3.1.8	Implementation on LLM Evaluation Metrics	9
56	3.1.9	Model Evaluation and Analysis of Results	9
57	3.1.10	Documentation	10
58	3.2	Calendar of Activities	10
59	4	Preliminary Results/System Prototype	11
60		References	14
61	A	Figures and Tables	17

62 List of Figures

<small>63</small>	4.1 Code snippet of the prototype	12
<small>64</small>	A.1 Code snippet of the prototype	18

65 List of Tables

66	3.1	Timetable of Activities	10
67	4.1	Difference between Generated Answers of the Base and LoRA Fine-	
68		Tuned Model	13
69	A.1	Difference between Generated Answers of the Base and LoRA Fine-	
70		Tuned Model	19

Chapter 1

Introduction

1.1 Overview

Language is how humans communicate and express themselves (Crystal & Robins, 2024). It is dynamic because there are endless structural possibilities, changes in word meanings, and new words created (Libretexts, 2021). Slang is a great example of the dynamic nature of language. Slang is an informal language used by people in the same social group (Fernández-Toro, 2016). It serves social purposes: to identify a group's members, communicate informally, and oppose established authority (McArthur, 2003). Slang is highly contextual and pervasive, even in non-standard English. (Roth-Gordon, 2020) Its figurative nature and how it twists the definitions of the words used in it make it hard for outsiders to understand (Mattiello, 2005).

In recent years, the internet has become a significant medium for the evolution and spread of language, giving rise to 'internet slang' (J. Liu, Zhang, & Li, 2023). Internet slang is a collection of everyday language forms used by diverse groups online (Barseghyan, 2014). Ujang et al. (2018, as cited in (binti Sabri, bin Hamdan, Nadarajan, & Shing, 2020)) state that Internet slang is not easily understood by people outside the social group or people who are not fluent in the language where slang is used. This phenomenon is particularly prominent among the younger generation (Maulidiya, Wijaya, Mauren, Adha, & Pandin, 2021), where they use it to communicate and interact with friends.

Today, Generation Alpha is the youngest generation. Generation Alpha refers to people born between 2010 and 2025. They were born into an era of rapid technological advancement, where digital devices and the internet are integral to

96 their daily lives (McCrindle & Fell, 2020). Generation Alpha is also called the
97 first true digital natives (Jukić & Škojo, 2021). They are expected to be the most
98 “technologically” skilled and most educated generation as they are the native
99 speakers of the language of the Internet (Prensky, 2001). According to the study
100 *Understanding Generation Alpha*, Generation Alpha is socially driven, which may
101 let them grow up to be creative and unconventional, potentially shaping them to
102 be assets in the future (Jha, 2020).

103 Since Generation Alpha was born with technology, the usage of Internet slang
104 has been prominent in this generation. However, it can create communication
105 barriers between older and younger generations (Venter, 2017 as cited in (Ghazali
106 & Abdullah, 2021)). The communication barriers caused by the usage of Inter-
107 net slang also affect people from the younger generation, especially individuals
108 who are less active on social media and have less exposure to them (Vacalares,
109 Salas, Babac, Cagalawan, & Calimpong, 2023). This gap highlights the need for
110 a tool that can bridge the generational divide, making it easier for individuals
111 to understand the language of Generation Alpha. By fostering a mutual under-
112 standing, such tool can promote more effective and harmonious interactions across
113 generations, enhancing relationships and reducing miscommunication.

114 1.2 Problem Statement

115 Internet slang fosters informal, relatable communication within the younger gen-
116 eration (Ghazali & Abdullah, 2021), especially Generation Alpha, but it presents
117 challenges in understanding for people outside this demographic. The gap in com-
118 prehension with older generations widens as internet slang evolves, often leading
119 to miscommunication affecting social relationships that contribute to the genera-
120 tional divide (Vacalares et al., 2023). A more specific translation tool developed
121 using language models can be used to bridge this divide.

122 By leveraging the ability of LLM to generate a more nuanced and properly
123 constructed answer, a better tool can be made to translate the slangs into proper
124 sentences. It has already been proven by the likes of GPT being modified and tai-
125 lored for use in several automated chatbots to provide customer service. However,
126 no such tool exists for slang translation of Generation Alpha, which arguably has
127 the most diverse slangs compared to other generations. The creation of this tool
128 will allow translating of such texts into formal sentences and help with bridging
129 the generational divide between them and older people, especially teachers.

130 1.3 Research Objectives

131 1.3.1 General Objectives

132 This study aims to modify an existing Large Language Model (LLM) for use in
133 the translation of Generation Alpha internet slang used by Filipino children in
134 social media.

135 1.3.2 Specific Objectives

- 136 • To create a dataset of sentences containing Gen Alpha slang and its formal
137 translation
- 138 • To create a Low Rank Adaptation (LoRA) implementation for fine-tuning
139 an existing model
- 140 • To fine-tune an existing LLM to translate sentences containing Gen Alpha
141 slang into formal sentences
- 142 • To evaluate the performance of the trained model and compare it to the
143 based model using several performance metrics

144 1.4 Scope and Limitations of the Research

145 This study will focus on the usage of internet slang by Filipino Generation Alpha,
146 with an emphasis on English language since it is widely use on different digital
147 platforms such as social media.

148 1.5 Significance of the Research

149 The study contributes to understanding the evolving linguistic landscape shaped
150 by internet slang, especially as used by Generation Alpha. Insights gained from
151 this study may aid educators, parents, and communication professionals in bridg-
152 ing inter-generational communication gaps and fostering better understanding
153 across age groups.

Chapter 2

Review of Related Literature

2.1 Communication Gap between Generations

Language is dynamic in nature thus, constantly evolving over time. One example of this behavior is the development of internet slang. Internet slang is a result of language variation and is often regarded as informal (S. Liu, Gui, Zuo, & Dai, 2019). In the study, *The Use of Online Slang for Independent Learning in English Vocabulary* (Ambarsaru, Amrullah, & Nawawi, 2020), students used internet slang to express their feelings and emotions, and to align their communication style with their peers.

However, this development has its challenges. It is suggested that younger generation should use slang to communicate with each instead of older generations because it might cause confusion between them (Jeresano & Carretero, 2022).

This miscommunication is prominent between generations with differences in linguistic familiarity as Suslak (Suslak, 2009) argues that age influences language use, noting that language evolves across generations. Supporting this, a study by Teng and Joo (teng & Joo, 2023) found that the older a person is, the less likely they are to understand internet language.

The studies showed that using internet slang improves relationship between those who use it. However, using internet slang for inter-generational communication can be a hindrance to proper and effective communication (Gonzaga, Racal, & Estrada, n.d.).

176 2.2 Existing Studies

177 Khazeni et al. (Heydari, Albadvi, & Khazeni, 2024) used deep learning to create a
178 model for translating Persian slang text into formal ones. The researchers explored
179 the challenges of translating Persian slang into English within the context of
180 film subtitling, specifically focusing on the performance of three neural machine
181 translation (NMT) systems, namely Google Translate, Targoman, and Farazin.
182 The primary interest of the paper lies in the understanding how these NMTs
183 systems handle the complexities of slang translation. It was revealed that the
184 NMT systems often struggle to capture the nuances of slang, leading to unnatural
185 and inaccurate translations. Targoman performed best in naturalness, it fell short
186 of human translation quality. This implied the need for specialized algorithms or
187 training data suitable for slang, and potentially human post-editing, to achieve
188 accurate and culturally appropriate translations in this domain.

189 The study by Nocon et al. (Nocon, Kho, & Arroyo, 2018) explores translating
190 Filipino colloquialisms, such as Conyo and Datkilab, into standardized Filipino,
191 addressing comprehension barriers for non-familiar speakers. Two machine
192 translation (MT) approaches were evaluated: Tensorflow’s Sequence-to-Sequence
193 model using Recurrent Neural Networks (RNNs) and Moses’ Phrase-based Statistical
194 MT. Moses outperformed Tensorflow on test data due to its handling of
195 phrase combinations and unfamiliar words, while Tensorflow excelled on training
196 data, indicating potential with refinement and more training data. The research
197 underscores the need for robust datasets and highlights the strengths of phrase-
198 based statistical MT in tackling slang translation challenges.

199 Ibrahim and Mustafa (Abdulstar Ibrahim & Shareef Mustafa, 2023) developed
200 a system to translate slang into formal language, addressing challenges posed by
201 slang’s informality and variability. Using updated datasets of slang words, formal
202 equivalents, and contextual sentences, they fine-tuned pre-trained models from
203 Hugging Face’s Transformer library. While the T5-base model showed promise
204 during training, it performed poorly in testing. In contrast, the “facebook/bart-
205 base” model excelled, demonstrating high accuracy and low loss values. The
206 study highlights the importance of fine-tuning and updated datasets for effective
207 slang translation and emphasizes the potential of transformer models like
208 “facebook/bart-base” in bridging informal and formal language gaps.

2.3 LoRA for Fine Tuning

Low Rank Adaptation, or LoRA, is an efficient Parameter Efficient Fine Tuning (PEFT) method proposed by Hu et al (Hu et al., 2021). It can significantly decrease the required storage for training while producing comparable results and in some cases, even outperforming other adaptation methods. In addition, it has minimal chance of catastrophic forgetting as the original weights are not being tampered with, unlike other finetuning methods. These factors make it a suitable option for slang translation as a quick yet accurate solution. In a study conducted by Zhao et al. (Zhao et al., 2024), they determined that some LLMs using LoRA for fine tuning can outperform GPT-4, one of the most advanced LLM models currently. A study by Nguyen et al. (Nguyen, Wilson, & Dalins, 2023) used LoRA in fine tuning a pre-trained Llama 2 7B model for text classification of a dataset that contains slang. They were able to create a more accurate model compared to models by existing studies at that time.

2.4 Chapter Summary

This chapter shows how generational differences create communication gaps, especially due to internet slang. Younger people tend to use slang to express emotions and connect with friends, but this can confuse older generations who aren't as familiar with these terms. Research shows that as language changes over time, older people are generally less likely to understand the newest internet language. To bridge this gap, some recent studies have utilized machine learning to translate slang into more standard language. For instance, Khazeni et al. (Heydari et al., 2024) used deep learning to translate Persian slang, while Nocon et al. (Nocon et al., 2018) created a Filipino slang translator using statistical models. Moreover, Ibrahim and Mustafa (Abdulstar Ibrahim & Shareef Mustafa, 2023) fine-tuned pre-trained models to learn slang meanings. One of the promising techniques for this is Low Rank Adaptation (LoRA), which is a fine-tuning method that keeps the original model stable while using less storage. Studies by Zhao et al. (Zhao et al., 2024) and Nguyen et al. (Nguyen et al., 2023) show that LoRA models are not only efficient but can even outperform advanced models like GPT-4 when it comes to slang translation and text classification.

240 Chapter 3

241 Research Methodology

242 This chapter lists and discusses the specific steps and activities that will be per-
243 formed to accomplish the project. The discussion covers the activities from pre-
244 proposal to Final SP Writing.

245 3.1 Research Activities

246 3.1.1 Creation of the dataset

247 A dataset of sentences containing Generation Alpha slangs and its formal trans-
248 lation or an approximation of will be created. This will involve data scraping,
249 use of existing datasets, or any other suitable methods of obtaining data. This
250 dataset will be used for the training and evaluation of the model. To ensure it is a
251 high quality dataset, it will be manually checked for accuracy and grammatically
252 correctness. It will also be checked for any potential biases that may exist in the
253 dataset or the data collection process.

254 A complete dataset of sentences containing Generation Alpha slangs is ex-
255 pected at the end of this task.

256 3.1.2 Identification of potential LLM to be used

257 We will be reading upon existing LLM comparison studies to identify potential
258 LLMs to be used for this study. We will be primarily using studies that used

dataset containing slangs as they are the most similar to our required dataset. A good potential model is zephyr-7b-beta due to its popularity, more open license, and number of parameters. Having 7B parameters allows the training of models on a 16GB GPU with a 4-bit quantization.

A model to use should be determined at the end of this task.

3.1.3 Lookup on available GPU on demand services

Available computing power rental services will be looked up for this study. As LLM training are a resource-intensive process, it is important to ensure that the necessary computing power is available. However, this computing power requires expensive equipment that might not see usage after the project is completed. Thus, it has been decided that it is better to rent the computing power for the duration of the project. A report on available GPU on demand services will be created using market research and price to computing power ratio.

3.1.4 Study on LoRA implementation for LLM

A thorough study on the implementation of LoRA for fine-tuning will be done. This includes learning the necessary steps, logic behind the idea, and other necessary information necessary for implementation. For this step, reading upon guide materials regarding fine-tuning and LoRA as well as existing studies will be done. We will be primarily using the guide provided by HuggingFace as it is one of the largest repositories for prebuilt LLMs. In addition, they also provided guides for fine-tuning models for specific purposes and has model specific guides.

3.1.5 Preprocessing of data

The dataset used for the fine-tuning of the model will be cleaned up. This will require removal of non essential information such as email addresses, URLs, etc. This is to ensure that the model can focus on learning the patterns between the slang and its formal translation without being affected by noise.

A clean dataset ready for tokenization is expected at the end of this task.

286 **3.1.6 Prototype implementation of LoRA**

287 A prototype implementation of LoRA will be created using a less demanding
288 model. This is to avoid incurring costs from constantly retraining the model due
289 to bugs in the code. It will be also developed on the same language as the final
290 implementation to avoid any issues with the code translation. As it is a prototype,
291 it will be used to create a foundation for the complete implementation of LoRA.
292 It will ensure that during the final implementation, there will be no issues with
293 the code and the model can be fairly evaluated.

294 For this task, Google Colab will be used as a platform of choice due to the free
295 cloud computing resource and the use of Jupyter notebook. In addition, Python
296 will be used as the language of choice to the abundance of available libraries for
297 training LLMs.

298 **3.1.7 Implementation of LoRA on selected model**

299 A full implementation of LoRA will be done using the previously created prototype
300 as a basis. This step will mostly involve tweaking the parameters used to train
301 the selected model and fixing any hidden bugs in the generated results.

302 **3.1.8 Implementation on LLM Evaluation Metrics**

303 A set of evaluation metrics will be used to determine if the fine-tuned model will
304 perform better than the base model. These metrics will be taken from existing
305 studies on LoRA finetuning and slang translation. It will serve as the primary
306 measure in which LLMs are compared with from each other. For this purpose,
307 Recall-Oriented Understudy for Gisting Evaluation (ROUGE) will be used to
308 score the generated output compared to ground truth. The use of LLM as a judge
309 might also be considered to directly compare the results of the fine-tuned and the
310 base model.

311 **3.1.9 Model Evaluation and Analysis of Results**

312 The model obtained from previous steps will be evaluated using the evaluation
313 metrics determined from the previous step. To do this, the testing set split of
314 the dataset will be used as the basis of evaluation. In addition, descriptive infor-

315 mation such as loss function per epoch and perplexity will be determined. This
 316 information will be used as supplement to evaluation metrics to determine if the
 317 fine-tuned model performed better than the base model.

318 3.1.10 Documentation

319 All members are tasked to provide accurate and detailed logs of their activities.
 320 This includes steps on the task they are working on, the status of the work being
 321 done, and the time spent on the task. It will serve both as documentation and as
 322 a progress tracker to determine how far the project is from being done. It will be
 323 done every week at the member's leisure.

324 3.2 Calendar of Activities

325 Table 3.1 shows a Gantt chart of the activities. Each bullet represents approxi-
 326 mately one week worth of activity.

Table 3.1: Timetable of Activities

Activities (2024-2025)	Dec	Jan	Feb	Mar	Apr	May	Jun
Creation of the dataset	•						
Identification of potential LLM to be used	•						
Lookup on available GPU on demand services	•						
Study on LoRA implementation for LLM	•						
Preprocessing of data	•••						
Prototype implementation of LoRA	•	••••					
Implementation of LoRA on selected model			••				
Implementation on LLM Evaluation Metrics			••				
Model Evaluation and Analysis of Results				••••			
Documentation	••	••••	••••	••••	••••		

327 Chapter 4

328 Preliminary Results/System 329 Prototype

330 A prototype LoRA implementation was created on Google Colab. This uses
331 zephyr-7b-beta model as the base model for finetuning and a part of the ultrachat
332 dataset as the training dataset. The zephyr-7b-beta model was chosen as it is one
333 of the best performing model after LoRA finetuning (Zhao et al., 2024).

334 This snippet is based on the fine-tuning guides available on HuggingFace and
335 the PEFT guides on several websites. This prototype only uses 100 iterations to
336 train the model based on the ultrachat dataset. Most of the parameters used in
337 this prototype is the same with the examples but it establishes a proof of concept
338 on the usage of Google Colab as a potential training platform and the ability of
339 LoRA to improve generative results.

Figure 4.1: Code snippet of the prototype

```
from transformers import TrainingArguments

YOUR_HF_USERNAME = "chichiboo"

output_dir = f"{YOUR_HF_USERNAME}/zephyr-7b-qlora-ultrachat"
per_device_train_batch_size = 4
gradient_accumulation_steps = 4
optim = "paged_adamw_32bit"
save_steps = 10
logging_steps = 10
learning_rate = 2e-4
max_grad_norm = 0.3
max_steps = 100
warmup_ratio = 0.03
lr_scheduler_type = "constant"

training_arguments = TrainingArguments(
    output_dir=output_dir,
    per_device_train_batch_size=per_device_train_batch_size,
    gradient_accumulation_steps=gradient_accumulation_steps,
    optim=optim,
    save_steps=save_steps,
    logging_steps=logging_steps,
    learning_rate=learning_rate,
    max_grad_norm=max_grad_norm,
    max_steps=max_steps,
    warmup_ratio=warmup_ratio,
    lr_scheduler_type=lr_scheduler_type,
```

Table 4.1: Difference between Generated Answers of the Base and LoRA Fine-Tuned Model

Before LoRA:

USER:	Can you explain contrastive learning in machine learning in simple terms for someone new to the field of ML?
ASSISTANT:	Contrastive learning is a machine learning technique that involves training a model to learn from data that is similar to the target data. The model is trained to identify patterns in the data that are similar to the target data, and to use those patterns to make predictions about new data.

After LoRA:

USER:	Can you explain contrastive learning in machine learning in simple terms for someone new to the field of ML?
ASSISTANT:	Contrastive learning is a technique in machine learning that uses a technique called contrastive loss to train a model. The idea is to train a model to distinguish between similar and dissimilar data. In contrastive learning, a model is trained to distinguish between a pair of similar data, such as a pair of images, and a pair of dissimilar data, such as a pair of different images. The model is trained to minimize the loss, which is the difference between the similarity and dissimilarity of the data. In contrast, traditional machine learning techniques, such as supervised learning, are trained to distinguish between different types of data, such as a classification or regression task. In summary, contrastive learning is a technique in machine learning that uses a technique called contrastive loss to train a model to distinguish between similar and dissimilar data, which is different from traditional machine learning techniques, such as supervised learning.

References

- Abdulstar Ibrahim, A., & Shareef Mustafa, B. (2023, Oct). Intelligent system to transformer slang words into formal words. *NTU Journal of Engineering and Technology*, 2(2). doi: 10.56286/ntujet.v2i2.689
- Ambarsaru, S., Amrullah, A., & Nawawi, N. (2020, Aug). The use of online slang for independent learning in english vocabulary. *Proceedings of the 1st Annual Conference on Education and Social Sciences (ACCESS 2019)*, 465, 295–297. doi: 10.2991/assehr.k.200827.074
- Barseghyan, L. (2014). *On some aspects of internet slang*. Retrieved from <https://api.semanticscholar.org/CorpusID:51730779>
- binti Sabri, N. A., bin Hamdan, S., Nadarajan, N.-T. M., & Shing, S. R. (2020, Jun). The usage of english internet slang among malaysians in social media. *Selangor Humaniora Review*, 4(1), 16-17.
- Crystal, D., & Robins, R. H. (2024, Oct). *Language*. Encyclopædia Britannica, inc. Retrieved from <https://www.britannica.com/topic/language>
- Fernández-Toro, M. (2016, Jun). *Exploring languages and cultures*. Retrieved from <https://www.open.edu/openlearn/languages/exploring-languages-and-cultures/content-section-3.2>
- Ghazali, N. M., & Abdullah, N. N. (2021, Dec). Slang language use in social media among malaysian youths: A sociolinguistic perspective. *International Young Scholars Journal of Languages*, 4(2), 69. Retrieved from https://www.iium.edu.my/media/77652/Slang%20Language%20Use%20in%20Social%20Media%20Among%20Malaysian%20Youths_A%20Sociolinguistic%20Perspective.pdf
- Gonzaga, L. M. T., Racal, J. T. Q., & Estrada, K. O. (n.d.). “forda convo ang ferson”: Analysis of gen z slang in the lens of batstateu faculty members.
- Heydari, M., Albadvi, A., & Khazeni, M. (2024). Persian slang text conversion to formal and deep learning of persian short texts on social media for sentiment classification. *Journal of Electrical and Computer Engineering Innovations (JECEI)*. Retrieved from https://jecei.sru.ac.ir/article_2172.html doi: 10.22061/jecei.2024.10745.731
- Hu, E. J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., ... Chen, W.

(2021). *Lora: Low-rank adaptation of large language models*. Retrieved from <https://arxiv.org/abs/2106.09685>

Jeresano, E., & Carretero, M. (2022, Feb). Digital culture and social media slang of gen z. *United International Journal for Research Technology*, 3(4), 11–25. doi: <http://dx.doi.org/10.1314/RG.2.2.36361.93285>

Jha, A. (2020, Jun). *Understanding generation alpha*. doi: 10.31219/osf.io/d2e8g

Jukić, R., & Škojo, T. (2021). The educational needs of the alpha generation. In *2021 44th international convention on information, communication and electronic technology (mipro)* (p. 564-569). doi: 10.23919/MIPRO52101.2021.9597106

Libretexts. (2021, Jul). *3.1.2: Functions of language*. Author. Retrieved from [https://socialsci.libretexts.org/Courses/American_River_College/SPEECH_361%3A_The_Communication_Experience_\(Coleman\)/03%3A_Verbal_Codes/3.01%3A_Verbal_Communication/3.1.02%3A_Functions_of_Language](https://socialsci.libretexts.org/Courses/American_River_College/SPEECH_361%3A_The_Communication_Experience_(Coleman)/03%3A_Verbal_Codes/3.01%3A_Verbal_Communication/3.1.02%3A_Functions_of_Language)

Liu, J., Zhang, X., & Li, H. (2023, Aug). Analysis of language phenomena in internet slang: A case study of internet dirty language. *Open Access Library Journal*, 10(08), 1–12. doi: 10.4236/oalib.1110484

Liu, S., Gui, D.-Y., Zuo, Y., & Dai, Y. (2019, Jun). Good slang or bad slang? embedding internet slang in persuasive advertising. *Frontiers in Psychology*, 10. doi: 10.3389/fpsyg.2019.01251

Mattiello, E. (2005). The pervasiveness of slang in standard and non-standard english.. Retrieved from <https://api.semanticscholar.org/CorpusID:140842571>

Maulidiya, R., Wijaya, S. E., Mauren, C., Adha, T. P., & Pandin, M. G. R. (2021, Dec). *Language development of slang in the younger generation in the digital era*. OSF Preprints. Retrieved from osf.io/xs7kd doi: 10.31219/osf.io/xs7kd

McArthur, T. (2003). *Concise oxford companion to the english language* (1st ed.). Oxford University Press.

McCrindle, M., & Fell, A. (2020). *Understanding generation alpha*. McCrindle Research Pty Ltd.

Nguyen, T. T., Wilson, C., & Dalins, J. (2023). *Fine-tuning llama 2 large language models for detecting online sexual predatory chats and abusive texts*. Retrieved from <https://arxiv.org/abs/2308.14683>

Nocon, N., Kho, N. M., & Arroyo, J. (2018, Oct). Building a filipino colloquialism translator using sequence-to-sequence model. *TENCON 2018 - 2018 IEEE Region 10 Conference*, 2199–2204. doi: 10.1109/tencon.2018.8650118

Prensky, M. (2001, Oct). Digital natives, digital immigrants. *On the Horizon*, 9(5). doi: <https://doi.org/10.1108/10748120110424816>

Roth-Gordon, J. (2020). Language and creativity: Slang. In *The international encyclopedia of linguistic anthropology* (p. 1-8). John Wi-

414 ley Sons, Ltd. Retrieved from [https://onlinelibrary.wiley.com/doi/](https://onlinelibrary.wiley.com/doi/abs/10.1002/9781118786093.iela0192)
 415 abs/10.1002/9781118786093.iela0192 doi: [https://doi.org/10.1002/](https://doi.org/10.1002/9781118786093.iela0192)
 416 9781118786093.iela0192
 417 Suslak, D. F. (2009). The sociolinguistic problem of generations. *Language Com-*
 418 *munication*, 29(3), 199-209. Retrieved from [https://www.sciencedirect](https://www.sciencedirect.com/science/article/pii/S0271530909000196)
 419 .com/science/article/pii/S0271530909000196 (Reflecting on language
 420 and culture fieldwork in the early 21st century) doi: [https://doi.org/](https://doi.org/10.1016/j.langcom.2009.02.003)
 421 10.1016/j.langcom.2009.02.003
 422 teng, C. E., & Joo, T. M. (2023). Is internet language a destroyer to communica-
 423 tion? In X.-S. Yang, R. S. Sherratt, N. Dey, & A. Joshi (Eds.), *Proceedings of*
 424 *eighth international congress on information and communication technology*
 425 (pp. 527–536). Singapore: Springer Nature Singapore.
 426 Vacalares, S. T., Salas, A. F. R., Babac, B. J. S., Cagalawan, A. L., & Calimpong,
 427 C. D. (2023, Jun). The intelligibility of internet slangs between millennials
 428 and gen zers: A comparative study. *International Journal of Science and*
 429 *Research Archive*, 9(1), 400–409. doi: 10.30574/ijrsra.2023.9.1.0456
 430 Zhao, J., Wang, T., Abid, W., Angus, G., Garg, A., Kinnison, J., ... Rishi, D.
 431 (2024). *Lora land: 310 fine-tuned llms that rival gpt-4, a technical report*.
 432 Retrieved from <https://arxiv.org/abs/2405.00732>

433 **Appendix A**

434 **Figures and Tables**

Figure A.1: Code snippet of the prototype

```
from transformers import TrainingArguments

YOUR_HF_USERNAME = "chichiboo"

output_dir = f"{YOUR_HF_USERNAME}/zephyr-7b-qlora-ultrachat"
per_device_train_batch_size = 4
gradient_accumulation_steps = 4
optim = "paged_adamw_32bit"
save_steps = 10
logging_steps = 10
learning_rate = 2e-4
max_grad_norm = 0.3
max_steps = 100
warmup_ratio = 0.03
lr_scheduler_type = "constant"

training_arguments = TrainingArguments(
    output_dir=output_dir,
    per_device_train_batch_size=per_device_train_batch_size,
    gradient_accumulation_steps=gradient_accumulation_steps,
    optim=optim,
    save_steps=save_steps,
    logging_steps=logging_steps,
    learning_rate=learning_rate,
    max_grad_norm=max_grad_norm,
    max_steps=max_steps,
    warmup_ratio=warmup_ratio,
    lr_scheduler_type=lr_scheduler_type,
```


Table A.1: Difference between Generated Answers of the Base and LoRA Fine-Tuned Model

Before LoRA:

USER:	Can you explain contrastive learning in machine learning in simple terms for someone new to the field of ML?
ASSISTANT:	Contrastive learning is a machine learning technique that involves training a model to learn from data that is similar to the target data. The model is trained to identify patterns in the data that are similar to the target data, and to use those patterns to make predictions about new data.

After LoRA:

USER:	Can you explain contrastive learning in machine learning in simple terms for someone new to the field of ML?
ASSISTANT:	Contrastive learning is a technique in machine learning that uses a technique called contrastive loss to train a model. The idea is to train a model to distinguish between similar and dissimilar data. In contrastive learning, a model is trained to distinguish between a pair of similar data, such as a pair of images, and a pair of dissimilar data, such as a pair of different images. The model is trained to minimize the loss, which is the difference between the similarity and dissimilarity of the data. In contrast, traditional machine learning techniques, such as supervised learning, are trained to distinguish between different types of data, such as a classification or regression task. In summary, contrastive learning is a technique in machine learning that uses a technique called contrastive loss to train a model to distinguish between similar and dissimilar data, which is different from traditional machine learning techniques, such as supervised learning.
