

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 From 150 to 200 words of short, direct and complete sentences, the abstract should
20 be informative enough to serve as a substitute for reading the entire SP document
21 itself. It states the rationale and the objectives of the research. In the final Special
22 Problem document (i.e., the document you'll submit for your final defense), the
23 abstract should also contain a description of your research results, findings, and
24 contribution(s).

25 Suggested keywords based on ACM Computing Classification system can be
26 found at https://dl.acm.org/ccs/ccs_flat.cfm

27 **Keywords:** Keyword 1, keyword 2, keyword 3, keyword 4, etc.

Contents

29	1 Introduction	1
30	1.1 Overview	1
31	1.2 Problem Statement	2
32	1.3 Research Objectives	2
33	1.3.1 General Objectives	2
34	1.3.2 Specific Objectives	3
35	1.4 Scope and Limitations of the Research	3
36	1.5 Significance of the Research	3
37	2 Review of Related Literature	4
38	2.1 Communication Gap between Generations	4
39	2.2 Existing Studies	4
40	2.3 LoRA for Fine Tuning	5
41	2.4 Chapter Summary	5
42	3 Research Methodology	7
43	3.1 Research Activities	7
44	3.1.1 Creation of the dataset	7

45	3.1.2	Identification of potential LLM to be used	7
46	3.1.3	Lookup on available GPU on demand services	8
47	3.1.4	Study on LoRA implementation for LLM	8
48	3.1.5	Preprocessing of data	8
49	3.1.6	Prototype implementation of LoRA	8
50	3.1.7	Implementation of LoRA on selected model	9
51	3.1.8	Implementation on LLM Evaluation Metrics	9
52	3.1.9	Model Evaluation and Analysis of Results	9
53	3.1.10	Documentation	9
54	3.2	Calendar of Activities	9
55	4	Preliminary Results/System Prototype	11
56		References	12
57	A	Appendix Title	15

58 List of Figures

59 List of Tables

<small>60</small>	3.1 Timetable of Activities	10
-------------------	---------------------------------------	----

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

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

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

104 1.2 Problem Statement

105 Internet slang fosters informal, relatable communication within the younger gen-
106 eration (Ghazali & Abdullah, 2021), especially Generation Alpha, but it presents
107 challenges in understanding for people outside this demographic. The gap in com-
108 prehension with older generations widens as internet slang evolves, often leading
109 to miscommunication affecting social relationships that contribute to the genera-
110 tional divide (Vacalares et al., 2023). A more specific translation tool developed
111 using language models use in many digital platforms can be used to bridge this
112 divide.

113 1.3 Research Objectives

114 1.3.1 General Objectives

115 This study aims to modify an existing Large Language Model (LLM) for use in
116 the translation of Generation Alpha internet slang used by Filipino children in

117 social media.

118 1.3.2 Specific Objectives

- 119 • To create a dataset of sentences containing Gen Alpha slang and its formal
120 translation
- 121 • To create a Low Rank Adaptation (LoRA) implementation for fine-tuning
122 an existing model
- 123 • To fine-tune an existing LLM to translate sentences containing Gen Alpha
124 slang into formal sentences
- 125 • To evaluate the performance of the trained model and compare it to the
126 based model using several performance metrics

127 1.4 Scope and Limitations of the Research

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

131 1.5 Significance of the Research

132 The study contributes to understanding the evolving linguistic landscape shaped
133 by internet slang, especially as used by Generation Alpha. Insights gained from
134 this study may aid educators, parents, and communication professionals in bridg-
135 ing inter-generational communication gaps and fostering better understanding
136 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 follow the communication style of their peers. However, 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. 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.

2.2 Existing Studies

Khazeni et al. used deep learning to create a model for translating Persian slang text into formal ones (Heydari, Albadvi, & Khazeni, 2024). They were able to create a model to convert texts from social media into sentiments for classification. Nocon et al. (Nocon, Kho, & Arroyo, 2018) created a Filipino colloquialism translator using Tensorflow's sequence-to-sequence model and Moses' phrase-based sta-

159 tistical machine translation. They found that the Moses model was able to create
160 a natural sounding translation, while the Tensorflow model often produced bad
161 sentences.

162 A slang translation system developed by Ibrahim and Mustafa (Abdulstar Ibrahim
163 & Shareef Mustafa, 2023) used models obtained from Hugging Face, a repository
164 of pre-trained models, and retrained it using a dataset containing slang and their
165 corresponding definition and example. They determined that these models can
166 be tweaked into learning the relationship between the slang and its meaning.

167 **2.3 LoRA for Fine Tuning**

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

181 **2.4 Chapter Summary**

182 This chapter shows how generational differences create communication gaps, espe-
183 cially due to internet slang. Younger people tend to use slang to express emotions
184 and connect with friends, but this can confuse older generations who aren't as
185 familiar with these terms. Research shows that as language changes over time,
186 older people are generally less likely to understand the newest internet language.
187 To bridge this gap, some recent studies have utilized machine learning to translate
188 slang into more standard language. For instance, Khazeni et al. (Heydari et al.,
189 2024) used deep learning to translate Persian slang, while Nocon et al. (Nocon et
190 al., 2018) created a Filipino slang translator using statistical models. Moreover,
191 Ibrahim and Mustafa (Abdulstar Ibrahim & Shareef Mustafa, 2023) fine-tuned

192 pre-trained models to learn slang meanings. One of the promising techniques for
193 this is Low Rank Adaptation (LoRA), which is a fine-tuning method that keeps
194 the original model stable while using less storage. Studies by Zhao et al. (Zhao
195 et al., 2024) and Nguyen et al. (Nguyen et al., 2023) show that LoRA models are
196 not only efficient but can even outperform advanced models like GPT-4 when it
197 comes to slang translation and text classification.

Chapter 3

Research Methodology

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

3.1 Research Activities

3.1.1 Creation of the dataset

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

3.1.2 Identification of potential LLM to be used

We will be reading upon existing LLM comparison studies to identify potential 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.

216 **3.1.3 Lookup on available GPU on demand services**

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

224 **3.1.4 Study on LoRA implementation for LLM**

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

232 **3.1.5 Preprocessing of data**

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

237 **3.1.6 Prototype implementation of LoRA**

238 A prototype implementation of LoRA will be created using a less demanding
239 model. This is to avoid incurring costs from constantly retraining the model due
240 to bugs in the code. It will be also developed on the same platform as the final
241 implementation to avoid any issues with the code running on different platforms.
242 As it is a prototype, it will be used to create a foundation for the complete
243 implementation of LoRA. It will ensure that during the final implementation,
244 there will be no issues with the code and the model can be fairly evaluated.

245 **3.1.7 Implementation of LoRA on selected model**

246 A full implementation of LoRA will be done using the previously created prototype
247 as a basis. Since it has been proven to work, this step will mostly involve fine-
248 tuning the selected model and fixing any hidden bugs.

249 **3.1.8 Implementation on LLM Evaluation Metrics**

250 A set of evaluation metrics will be used to determine if the fine-tuned model will
251 perform better than the base model. These metrics will be taken from existing
252 studies on LoRA finetuning and slang translation. It will serve as the primary
253 measure in which LLMs are compared with from each other.

254 **3.1.9 Model Evaluation and Analysis of Results**

255 The model obtained from previous steps will be evaluated using the evaluation
256 metrics determined from the previous step. To do this, the testing set split of the
257 dataset will be used as the basis of evaluation. In addition, descriptive information
258 such as loss function per epoch, accuracy, precision, recall, and F1 score will be
259 determined. This information will be used as supplement to evaluation metrics to
260 determine if the fine-tuned model will perform better than the base model.

261 **3.1.10 Documentation**

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

267 **3.2 Calendar of Activities**

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

Table 3.1: Timetable of Activities

Activities (2024-2025)	Nov	Dec	Jan	Feb	Mar	Apr	May
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	••	••••	••••	••••	••••		

270 Chapter 4

271 Preliminary Results/System 272 Prototype

273 This chapter presents the preliminary results or the system prototype of your SP.
274 Include screenshots, tables, or graphs and provide the discussion of results.

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
- 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 Wiley Sons, Ltd. Retrieved from <https://onlinelibrary.wiley.com/doi/abs/10.1002/9781118786093.iela0192> doi: <https://doi.org/10.1002/>

349 9781118786093.iela0192
 350 Suslak, D. F. (2009). The sociolinguistic problem of generations. *Language Com-*
 351 *munication*, 29(3), 199-209. Retrieved from [https://www.sciencedirect](https://www.sciencedirect.com/science/article/pii/S0271530909000196)
 352 [.com/science/article/pii/S0271530909000196](https://www.sciencedirect.com/science/article/pii/S0271530909000196) (Reflecting on language
 353 and culture fieldwork in the early 21st century) doi: [https://doi.org/](https://doi.org/10.1016/j.langcom.2009.02.003)
 354 [10.1016/j.langcom.2009.02.003](https://doi.org/10.1016/j.langcom.2009.02.003)
 355 teng, C. E., & Joo, T. M. (2023). Is internet language a destroyer to communica-
 356 tion? In X.-S. Yang, R. S. Sherratt, N. Dey, & A. Joshi (Eds.), *Proceedings of*
 357 *eighth international congress on information and communication technology*
 358 (pp. 527–536). Singapore: Springer Nature Singapore.
 359 Vacalares, S. T., Salas, A. F. R., Babac, B. J. S., Cagalawan, A. L., & Calimpong,
 360 C. D. (2023, Jun). The intelligibility of internet slangs between millennials
 361 and gen zers: A comparative study. *International Journal of Science and*
 362 *Research Archive*, 9(1), 400–409. doi: 10.30574/ijrsra.2023.9.1.0456
 363 Zhao, J., Wang, T., Abid, W., Angus, G., Garg, A., Kinnison, J., ... Rishi, D.
 364 (2024). *Lora land: 310 fine-tuned llms that rival gpt-4, a technical report*.
 365 Retrieved from <https://arxiv.org/abs/2405.00732>

³⁶⁶ **Appendix A**

³⁶⁷ **Appendix Title**