

1
2

LOST IN TRANSLATION: UNDERSTANDING GENERATION ALPHA INTERNET SLANG

3
4
5
6
7
8

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

9
10
11

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

12
13
14

FLAUTA, Neil Bryan
GIMENO, Ashley Joy
GIMENO, Carl Jorenz

15
16

Francis DIMZON
Adviser

17

November 6, 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	1.4 Specific Objectives	3
35	1.5 Scope and Limitations of the Research	3
36	1.6 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	Testing and Analysis of Results	9
53	3.1.10	Documentation	9
54	3.2	Calendar of Activities	10
55	References		11

56 List of Figures

57 List of Tables

<small>58</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. Its figurative nature and how it twists the definitions of the words used in it make it hard for outsiders to understand.

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 their daily lives (McCrindle & Fell, 2020). Generation Alpha is also called the

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

90 Since Generation Alpha was born with technology, the usage of Internet slang
91 has been prominent in this generation. However, it can create communication bar-
92 riers between older and younger generations (Venter, 2017 as cited in (Ghazali &
93 Abdullah, 2021)). A study by Vargas and Barbella (Vargas & Marbella, 2023) in-
94 vestigated Generation Alpha's Filipino vocabulary and found that it often creates
95 misunderstandings for students and teachers—who are less familiar with internet
96 slang.

97 1.2 Problem Statement

98 Internet slang fosters informal, relatable communication within the younger gen-
99 eration (Ghazali & Abdullah, 2021), especially Generation Alpha, but it presents
100 challenges in understanding for people outside this demographic. The gap in com-
101 prehension with older generations widens as internet slang evolves, often leading
102 to miscommunication affecting social relationships that contribute to the genera-
103 tional divide (?, ?). This study investigates the communication barriers internet
104 slang creates, particularly between Generation Alpha and older generations, and
105 explores possible solutions to bridge this gap.

106 1.3 Research Objectives

107 1.3.1 General Objectives

108 This study aims to modify an existing LLM for use in the translation of Generation
109 Alpha internet slang used by Filipino children in social media.

110 1.4 Specific Objectives

- 111 • To create a dataset of sentences containing gen alpha slang and its formal
112 translation
- 113 • To create a LoRA implementation for fine-tuning an existing model
- 114 • To fine-tune an existing LLM to translate sentences containing gen alpha
115 slang into formal sentences
- 116 • To evaluate the performance of the trained model and compare it to the
117 based model using several performance metrics

118 1.5 Scope and Limitations of the Research

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

122 1.6 Significance of the Research

123 The study contributes to understanding the evolving linguistic landscape shaped
124 by internet slang, especially as used by Generation Alpha. Insights gained from
125 this study may aid educators, parents, and communication professionals in bridg-
126 ing intergenerational communication gaps and fostering better understanding across
127 age groups.

Chapter 2

Review of Related Literature

2.1 Communication Gap between Generations

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* (Ambarsari, Amrullah, & Nawawi, 2020), students used internet slang to express their feelings and emotions and because their friends also use it, However, it suggests 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 statistical machine translation. They found that the Moses model was able to create a natural sounding translation, while the Tensorflow model often produced bad

150 sentences.

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

156 2.3 LoRa for Fine Tuning

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

170 2.4 Chapter Summary

171 This chapter shows how generational differences create communication gaps, espe-
172 cially due to internet slang. Younger people tend to use slang to express emotions
173 and connect with friends, but this can confuse older generations who aren't as
174 familiar with these terms. Research shows that as language changes over time,
175 older people are generally less likely to understand the newest internet language.
176 To bridge this gap, some recent studies have utilized machine learning to translate
177 slang into more standard language. For instance, Khazeni et al. (Heydari et al.,
178 2024) used deep learning to translate Persian slang, while Nocon et al. (Nocon et
179 al., 2018) created a Filipino slang translator using statistical models. Moreover,
180 Ibrahim and Mustafa (Abdulstar Ibrahim & Shareef Mustafa, 2023) fine-tuned
181 pre-trained models to learn slang meanings. One of the promising techniques for
182 this is Low Rank Adaptation (LoRA), which is a fine-tuning method that keeps

183 the original model stable while using less storage. Studies by Zhao et al. (Zhao
184 et al., 2024) and Nguyen et al. (Nguyen et al., 2023) show that LoRA models are
185 not only efficient but can even outperform advanced models like GPT-4 when it
186 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

Ashley Joy Gimeno will be in-charge of creating a dataset of sentences containing Generation Alpha slangs and providing a formal translation of said sentence. This might involve data scraping ,reliance on existing dataset, or any other suitable method of obtaining it. This should last for a week and will serve as the training and testing information for the large language model during fine-tuning.

3.1.2 Identification of potential LLM to be used.

Carl Jorenz Gimeno will be tasked with finding potential models for the project and comparing them based on existing results. Having existing study using LoRA would be appreciated but does not solely determine it being used for this study. This should last for a week and a report on the prospect models will be created, detailing their strengths and weaknesses.

205 **3.1.3 Lookup on available GPU on demand services**

206 Neil Bryan Flauta will be tasked to find any reputable services that sell computing
207 power. This is essential as the group does not have direct access to hardware
208 necessary to fine-tune the selected model.

209 **3.1.4 Study on LoRA implementation for LLM**

210 Carl Jorenz Gimeno will be in-charge of studying on how LoRA is implemented
211 to LLMs. This will require reading various guides, primarily ones created by
212 HuggingFace as they are the creators of the model to be used and has several
213 in-depth guides in fine-tuning models in general. This should last a week and
214 Carl Jorenz Gimeno is expected to have the required knowledge by the end of it.

215 **3.1.5 Preprocessing of data**

216 Ashley Joy Gimeno will be tasked with preprocessing the data. Their task is to
217 ensure that all sentences contain at least one slang and all the formal translation
218 of the sentence is both grammatically correct and semantically correct. As LoRA
219 does not tamper with existing knowledge of the model (Hu et al., 2021), we are free
220 to focus on teaching the model the slang while leveraging its original knowledge
221 to provide proper sentences. In addition, after cleaning up the dataset, it will be
222 split into a training and testing set. This task should last 2-3 weeks or longer
223 based on the number of data points collected. A dataset ready for fine-tuning
224 should be available at the end

225 **3.1.6 Prototype implementation of LoRA**

226 Carl Jorenz Gimeno will be tasked with the implementation of LoRA on the
227 selected model. This includes applying a prototype to a smaller model and testing
228 the results. Carl Jorenz Gimeno may also opt to use qLoRA instead for the smaller
229 memory requirements at the cost of runtime (Raschka, 2023). Carl Jorenz Gimeno
230 must implement it using the selected computing service to prevent future changes
231 to adjust to the platform. This should last 4-5 weeks but could take more based
232 on the difficulty of actual implementation. It will serve as the basis of the proper
233 implementation of LoRA on the selected model to prevent longer testing with a

234 massive LLM. A working and correct implementation of LoRA should be available
235 at the end.

236 **3.1.7 Implementation of LoRA on selected model**

237 Neil Bryan Flauta will be tasked with the final implementation of LoRA on the
238 selected model, based on the prototype created. This should only last 1-2 weeks
239 because the code is already proven and tested as functional. A fine-tuned model
240 is expected to be complete at the end.

241 **3.1.8 Implementation on LLM Evaluation Metrics**

242 Neil Bryan Flauta will be tasked with studying the evaluation metrics used in
243 LLMs as well as create an implementation of such metrics. It will serve as a basis
244 in which we will compare the fine-tuned model with the base model. This should
245 take 2 weeks and a complete implementation of the metrics should be available at
246 the end.

247 **3.1.9 Testing and Analysis of Results**

248 Ashley Joy Gimeno will be tasked with testing the trained model using the testing
249 set on the dataset. This would include descriptive information regarding the model
250 and comparison with the original model.

251 **3.1.10 Documentation**

252 All members are tasked to provide accurate and detailed logs of their activities.
253 It will serve both as documentation and as a progress tracker to determine how
254 far the project is from being done. It will be done every week at the member's
255 leisure.

256 3.2 Calendar of Activities

257 Table 3.1 shows a Gantt chart of the activities. Each bullet represents approxi-
 258 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			••				
Testing and Analysis of Results				••••			
Documentation	••	••••	••••	••••	••••		

259 References

- 260 Abdulstar Ibrahim, A., & Shareef Mustafa, B. (2023, Oct). Intelligent system
261 to transform slang words into formal words. *NTU Journal of Engineering
262 and Technology*, 2(2). doi: 10.56286/ntujet.v2i2.689
- 263 Ambarsari, S., Amrullah, A., & Nawawi, N. (2020, Aug). The use of online
264 slang for independent learning in english vocabulary. *Proceedings of the 1st
265 Annual Conference on Education and Social Sciences (ACCESS 2019)*, 465,
266 295–297. doi: 10.2991/assehr.k.200827.074
- 267 Barseghyan, L. (2014). *On some aspects of internet slang*. Retrieved from
268 <https://api.semanticscholar.org/CorpusID:51730779>
- 269 binti Sabri, N. A., bin Hamdan, S., Nadarajan, N.-T. M., & Shing, S. R. (2020,
270 Jun). The usage of english internet slang among malaysians in social media.
271 *Selangor Humaniora Review*, 4(1), 16-17.
- 272 Crystal, D., & Robins, R. H. (2024, Oct). *Language*. Encyclopædia Britannica,
273 inc. Retrieved from <https://www.britannica.com/topic/language>
- 274 Fernández-Toro, M. (2016, Jun). *Exploring languages and cultures*. Re-
275 trieved from [https://www.open.edu/openlearn/languages/exploring
276 -languages-and-cultures/content-section-3.2](https://www.open.edu/openlearn/languages/exploring-languages-and-cultures/content-section-3.2)
- 277 Ghazali, N. M., & Abdullah, N. N. (2021, Dec). Slang language use
278 in social media among malaysian youths: A sociolinguistic per-
279 spective. *International Young Scholars Journal of Languages*,
280 4(2), 69. Retrieved from [https://www.iium.edu.my/media/
281 77652/Slang%20Language%20Use%20in%20Social%20Media%20Among%
282 20Malaysian%20Youths_A%20Sociolinguistic%20Perspective.pdf](https://www.iium.edu.my/media/77652/Slang%20Language%20Use%20in%20Social%20Media%20Among%20Malaysian%20Youths_A%20Sociolinguistic%20Perspective.pdf)
- 283 Heydari, M., Albadvi, A., & Khazeni, M. (2024). Persian slang text conversion to
284 formal and deep learning of persian short texts on social media for sentiment
285 classification. *Journal of Electrical and Computer Engineering Innovations
286 (JECEI)*. Retrieved from <https://jecei.sru.ac.ir/article.2172.html>
287 doi: 10.22061/jecei.2024.10745.731
- 288 Hu, E. J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., ... Chen, W.
289 (2021). *Lora: Low-rank adaptation of large language models*. Retrieved
290 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
- 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>
- Suslak, D. F. (2009). The sociolinguistic problem of generations. *Language Communication*, 29(3), 199-209. Retrieved from <https://www.sciencedirect.com/science/article/pii/S0271530909000196> (Reflecting on language and culture fieldwork in the early 21st century) doi: <https://doi.org/10.1016/j.langcom.2009.02.003>
- Teng, C. E., & Joo, T. M. (2023). Is internet language a destroyer to communication? In X.-S. Yang, R. S. Sherratt, N. Dey, & A. Joshi (Eds.), *Proceedings of*

333 *eighth international congress on information and communication technology*
 334 (pp. 527–536). Singapore: Springer Nature Singapore.
 335 Vargas, A., & Marbella, F. (2023, Sep). Bokabularyong generation alpha sa
 336 pakikipagtalastasang filipino. *International Journal of Research Studies in*
 337 *Education*, 12(8), 57–69. doi: <http://dx.doi.org/10.5861/ijrse.2023.62>
 338 Zhao, J., Wang, T., Abid, W., Angus, G., Garg, A., Kinnison, J., ... Rishi, D.
 339 (2024). *Lora land: 310 fine-tuned llms that rival gpt-4, a technical report*.
 340 Retrieved from <https://arxiv.org/abs/2405.00732>