

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 November 27, 2024

Contents

19	1 Introduction	1
20	1.1 Overview	1
21	1.2 Problem Statement	2
22	1.3 Research Objectives	2
23	1.3.1 General Objectives	2
24	1.4 Specific Objectives	3
25	1.5 Scope and Limitations of the Research	3
26	1.6 Significance of the Research	3
27	2 Review of Related Literature	4
28	2.1 Communication Gap between Generations	4
29	2.2 Existing Studies	4
30	2.3 LoRA for Fine Tuning	5
31	2.4 Chapter Summary	5
32	3 Research Methodology	7
33	3.1 Research Activities	7
34	3.1.1 Creation of the dataset	7

35	3.1.2	Identification of potential LLM to be used	7
36	3.1.3	Lookup on available GPU on demand services	8
37	3.1.4	Study on LoRA implementation for LLM	8
38	3.1.5	Preprocessing of data	8
39	3.1.6	Prototype implementation of LoRA	8
40	3.1.7	Implementation of LoRA on selected model	9
41	3.1.8	Implementation on LLM Evaluation Metrics	9
42	3.1.9	Model Evaluation and Analysis of Results	9
43	3.1.10	Documentation	9
44	3.2	Calendar of Activities	9

45 List of Tables

<small>46</small>	3.1 Timetable of Activities	10
-------------------	---------------------------------------	----

Chapter 1

Introduction

1.1 Overview

Language is how humans communicate and express themselves (?). It is dynamic because there are endless structural possibilities, changes in word meanings, and new words created (?). Slang is a great example of the dynamic nature of language. Slang is an informal language used by people in the same social group (?). It serves social purposes: to identify a group's members, communicate informally, and oppose established authority (?). 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' (?). Internet slang is a collection of everyday language forms used by diverse groups online (?). Ujang et al. (2018, as cited in (?)) 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 (?), 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 (?). Generation Alpha is also called the first true digital natives (?). They are expected to be the most "technologically" skilled and most educated generation as they are the native speakers of the language of the Internet (?). According to the study *Understanding Generation Alpha*, Generation Alpha

72 is socially driven, which may let them grow up to be creative and unconventional,
73 potentially shaping them to be assets in the future (?, ?).

74 Since Generation Alpha was born with technology, the usage of Internet slang
75 has been prominent in this generation. However, it can create communication bar-
76 riers between older and younger generations (Venter, 2017 as cited in (?, ?)). The
77 communication barriers caused by the usage of Internet slang also affect people
78 from the younger generation, especially individuals who are less active on social
79 media and have less exposure to them (?, ?). This gap highlights the need for a
80 tool that can bridge the generational divide, making it easier for individuals to
81 understand the language of Generation Alpha. By fostering a mutual understand-
82 ing, such a tool can promote more effective and harmonious interactions across
83 generations, enhancing relationships and reducing miscommunication.

84 1.2 Problem Statement

85 Internet slang fosters informal, relatable communication within the younger gen-
86 eration (?, ?), especially Generation Alpha, but it presents challenges in under-
87 standing for people outside this demographic. The gap in comprehension with
88 older generations widens as internet slang evolves, often leading to miscommuni-
89 cation affecting social relationships that contribute to the generational divide (?,
90 ?). A more specific translation tool developed using language models use in many
91 digital platforms can be used to bridge this divide.

92 1.3 Research Objectives

93 1.3.1 General Objectives

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

97 1.4 Specific Objectives

- 98 • To create a dataset of sentences containing gen alpha slang and its formal
99 translation
- 100 • To create a Low Rank Adaptation (LoRA) implementation for fine-tuning
101 an existing model
- 102 • To fine-tune an existing LLM to translate sentences containing gen alpha
103 slang into formal sentences
- 104 • To evaluate the performance of the trained model and compare it to the
105 based model using several performance metrics

106 1.5 Scope and Limitations of the Research

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

110 1.6 Significance of the Research

111 The study contributes to understanding the evolving linguistic landscape shaped
112 by internet slang, especially as used by Generation Alpha. Insights gained from
113 this study may aid educators, parents, and communication professionals in bridg-
114 ing intergenerational communication gaps and fostering better understanding across
115 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 (?). In the study, *The Use of Online Slang for Independent Learning in English Vocabulary* (?), 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 (?).

This miscommunication is prominent between generations. Suslak (?) argues that age influences language use, noting that language evolves across generations. Supporting this, a study by Teng and Joo (?) 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 (?). They were able to create a model to convert texts from social media into sentiments for classification. Nocon et al. (?) 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 sentences.

137 A slang translation system developed by Ibrahim and Mustafa (?, ?) used
138 models obtained from Hugging Face, a repository of pre-trained models, and re-
139 trained it using a dataset containing slang and their corresponding definition and
140 example. They determined that these models can be tweaked into learning the
141 relationship between the slang and its meaning.

142 2.3 LoRA for Fine Tuning

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

156 2.4 Chapter Summary

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

171 translation and text classification.

172 Chapter 3

173 Research Methodology

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

177 3.1 Research Activities

178 3.1.1 Creation of the dataset

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

186 3.1.2 Identification of potential LLM to be used

187 We will be reading upon existing LLM comparison studies to identify potential
188 LLMs to be used for this study. We will be primarily using studies that used
189 dataset containing slangs as they are the most similar to our required dataset.

190 **3.1.3 Lookup on available GPU on demand services**

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

198 **3.1.4 Study on LoRA implementation for LLM**

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

206 **3.1.5 Preprocessing of data**

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

211 **3.1.6 Prototype implementation of LoRA**

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

219 **3.1.7 Implementation of LoRA on selected model**

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

223 **3.1.8 Implementation on LLM Evaluation Metrics**

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

228 **3.1.9 Model Evaluation and Analysis of Results**

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

235 **3.1.10 Documentation**

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

241 **3.2 Calendar of Activities**

242 Table 3.1 shows a Gantt chart of the activities. Each bullet represents approxi-
243 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	••	••••	••••	••••	••••		