

# Korean "Airbnb" Script Deobfuscation Using Sequence-to-Sequence Models.

**Abdul Rafay**  
Boston University  
rafaya@bu.edu

**Mehmet Bora Sarioglu**  
Boston University  
Sarioglu@bu.edu

## Abstract

Korean online communities have developed a unique form of phonetic text coding for discreet communication. This project aims to develop an automated decoder for this coded Korean text. We will create a dataset of coded and decoded text, train a model to recognize transformation patterns, and deploy a functional tool for real-time decoding. This tool will bridge the communication gap, improving accessibility and understanding of online discourse.

## 1 Introduction

### 1.1 Problem Statement

Koreans online often use a phonetic text coding system popularly known as "airbnb script" to communicate candidly, sometimes obscuring messages. By subtly altering word phonetics, users convey meaning while maintaining a degree of privacy. This practice is common in online reviews, social media, and forums. While effective for native speakers, this coded language presents challenges for automated systems and non-Koreans.

### 1.2 Objective

Our goal for this project is to experiment and develop a decoder tool of the "airbnb script". The tool features a simple interface where users input encoded Korean text and the decoded standard Korean is displayed. The tool will provide real-time decoding with immediate feedback as users type or paste text.

### 1.3 Potential Impact

This project may have significant potential impact on facilitating understanding of Korean online discourse. It may also provide a valuable resource for studying Korean text phonetic transformations and at the core will offer a solution to a niche problem allowing for further applications. For businesses,

this tool will enable better understanding of customer feedback that may contain encoded Korean text.

### 1.4 About Text Obfuscation

Text obfuscation—the deliberate alteration of text to make it difficult for automated systems to process while remaining comprehensible to humans—poses significant challenges across multiple domains. Although we are interested in specific case of Korean "airbnb script", for better understanding of the concept, we can try to understand it through the context of English first. Common examples include replacing characters with visually similar ones (e.g., "noon" with "nouan"), substituting words with phonetically similar alternatives (e.g., "meet me later" → "meat mi letter"), or other transformations that preserve human readability while confusing automated systems.

These obfuscation techniques are frequently employed to:

- Evade content moderation systems
- Bypass plagiarism detection
- Circumvent keyword filtering
- Avoid automated content analysis

**Original:** Please meet me at the coffee shop tomorrow at noon.

**Obfuscated:** Pleezi meat mie it tha koffee shop timariw et noun.

Figure 1: Example of obfuscation.

Our research aims to develop a robust system for deobfuscation—converting obfuscated text back to its standard form—by leveraging phonetic representations as an intermediate step. This approach is motivated by the observation that many obfuscation techniques maintain phonetic similarity to preserve human readability.



### 2.4.3 Neural Machine Translation

Treating obfuscated text as a "source language" to be translated and treating the original text as the "target language" we trained character level and sub-character level LSTM models. They were significantly stronger than the statistical models but still came short of providing a useful degree of de-obfuscation.

## 3 Method

Our final approach combines phonetic preprocessing of the obfuscated text in order to mimic human thought process or behavior with a character level sequence-to-sequence model to recover original text from obfuscated versions.

### 3.1 Phonetic Preprocessing

How can Koreans understand encoded text, without non native speaker knowing? Answer to this was the pronunciation. Koreans automatically pronounce those encoded text in their head which happens to sound similar to regular phrases. From this realization, we tried to look for ways to convert encoded text into phonetic representations that could be easily parsed or trained.

#### 3.1.1 g2pK Module

Hangul, the main script for Korean, is phonetic, but the pronunciation rules are notoriously complicated. That's why we employed the g2pK (Grapheme-to-Phoneme) library specialized for Korean (available at <https://github.com/Kyubyong/g2pK>) to convert both original and obfuscated texts into phonetic representations that matches the actual sound of the phrases. This conversion maps visually distinct but phonetically similar characters and character sequences to the nearly similar phonetic representation, providing our model with a more consistent input space.

**Original:** 어제는 날씨가 맑았는데...  
**Obfuscated:** 엇째눈 날쭈카 맑앗눈데...  
**g2pK conversion:** 얼째눈 날쭈카 말간눈데...

Figure 3: Example of g2pK Conversion.

So, this would convert obfuscated text to a phonetic representation that captures their sound patterns, regardless of the specific characters used to create those sounds, making it more similar to original text.

### 3.1.2 Hangul Romanization

Another way of phonetic representation explored for potential integration into pipeline was romanization of Korean alphabets. There was a python module (available at <https://github.com/osori/korean-romanizer>) that comply with the rule developed by the National Institute of Korean Language, the official romanization system being used in the Republic of Korea. However, we ended up using g2pK only.

### 3.2 Character-Level Tokenization

Rather than using subword or word-level tokenization, we implement a character-level tokenizer that treats each character as a distinct token. Our tokenizer includes special tokens for padding (<pad>), unknown characters (<unk>), sequence start (<s>), and sequence end (</s>).

### 3.3 Model Architecture

We employ a BART-based (Lewis et al., 2020) sequence-to-sequence architecture with the following specifications:

Phonetic-Aware Text Deobfuscation Model Architecture

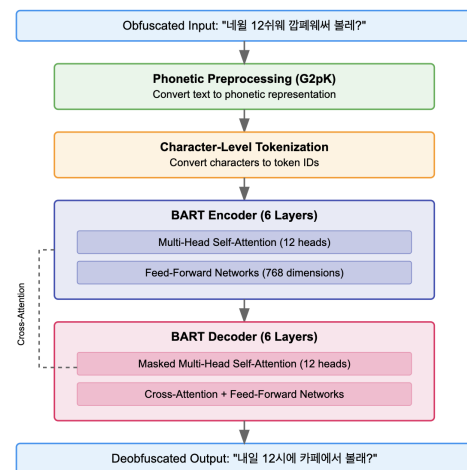


Figure 4: Model Architecture Diagram

- Vocabulary size based on unique characters in our dataset
- 6 encoder layers and 6 decoder layers
- 768-dimensional hidden representations
- 12 attention heads in both encoder and decoder

- Maximum position embeddings of 1024 tokens

The model inputs are the phonetic representations of obfuscated text, and it is trained to output the phonetic representations of the corresponding original text. In addition to this we have also tried to implement this pipeline using KoBART a pre-trained model specialized in the context of Korean language. However, we focused on base BART with our configurations to it. KoBART was promising and had better metrics compared to base BART with small data, however the lack of time and the inconsistency of the SCC cluster resulted in us not being able to train it on a large dataset, therefore we decided not to use it as our final model.

### 3.4 Training Procedure

We trained our model using cross-entropy loss with the following hyperparameters:

- Learning rate:  $3e-5$
- Batch size: 16 per device
- Weight decay: 0.01
- Training epochs: 10
- FP16 precision (when GPU is available)

Training was done on SCC, and took 14 hours to complete.

## 4 Experiments

### 4.1 Dataset

We created a large-scale dataset of paired original and obfuscated texts. For this we borrowed text data from **KoWiki** dataset, an Korean Wikipedia counterpart. It consists of the headers and text of Wikipedia in Korean. To generate the obfuscated dataset, we had to create the obfuscator ourselves, before we could train our model. We scraped the the most popular tool by hanmesoft (available at <https://airbnbify.hanmesoft.com/>) and retrieved the javascript file with the obfuscation logic that was being used in the frontend of the website and used the same code to write our own obfuscation code. The dataset contains approximately 10 million text pairs, of which we used fraction for training, with each pair consisting of an original text and its obfuscated version. The obfuscation includes character substitutions, phonetic spelling variations, and

Dataset	Lines	Size
Train	26,794,425	1.7G
Dev	130,419	7.7M
Test	134,340	8.4M

Table 1: KoWiki Dataset Information

other transformations that preserve the general phonetic structure.

Encoded Korean Texts	
original_text	encoded_text
세계수의 미국 시리즈에 전통	세계수의 미국 씨뤼츨엿 전통
2편 제왕의 성배부터 등장했	2편 제왕의 성배뽕뽕 등장했
세계수의 모험가들이 탐험하는	썰개수원 모험가들위 탐험한
그러나 분배할 수 있는 스킬	국렬날 뽕뽕할 슨 잇는 수퀼
다만 채집 시스템은 신세계수	달뽕 채집 신슐렘은 신세겻수
채집용 캐릭터들로 이루어진	썰침용 개릭터들로 위뤼어췌
필드 전투를 회피하면서 채집	필드 전투를 힘뽕함면췌 채췌
작품마다 !!아앗!!의 세세한 모	썰썰마다 !!아알!!위 췌세한 모
그 악랄함은 첫 등장한 작품이	고 악랄함은 첫 통췌한 작품이
게다가 이럴 때 쓰라고 있는 레	겐뽕가 위렐 뽕 췌락고 잇는
3편, 4편에는 숨통이 트이게	3편, 4편네는 슨뽕위 트뤼게
는 메시지가 뜨고 이때 운이 좋	는 뽕췌췌갈 뜨고 잇뽕 운잇
단 4편은 움직이지 않고 채집	뽕 4편는 움췌잉췌 안고 췌집
그리고 난이도 CASUAL로	클뤼고 난뽕뽕 CASUAL론
그나마 위험감지 먹통과 같은	구나마 위험감췌 먹뽕과 갓뽕

Figure 5: Example Pairs from Dataset

### 4.2 Evaluation Metrics

We evaluate our model using the following metrics:

**Character Error Rate (CER):** The Levenshtein distance between predicted and original texts, normalized by the length of the original text.

**Word Error Rate (WER):** The WER is derived from the Levenshtein distance, working at the word level instead of the phoneme level.

**BLEU Score:** To measure the general similarity between predicted and original texts.

### 4.3 Baselines

Before settling with our phonetic-aware approach, we tried following baselines:

- **LSTM:** We set up a character-level sequence to sequence where we created our of Encoder and Decoder architecture with LSTMs that we wrote ourself using pytorch.

- **Rule-based Normalization:** A system using handcrafted rules for common obfuscation patterns (e.g. "ㅈ" → "ㅉ", "ㅊ" → "ㅊㅊ").
- **Character-level JAMO Transformer:** We used character-level tokenization on decomposed Hangul characters (Jamos) to enhance phonetic representation. We also introduced Jamo embeddings to improve context understanding. Special tokens handle sequence padding and start/end markers.

#### 4.4 Results

Our phonetic-aware model outperforms all baselines, with a particularly significant improvement in BLEU score. We also did poll with native Korean speakers on our model’s performance. The feedback we got implied that the model was doing a good job in general, on a scale out of 10 we got an average around 8 in terms of deobfuscating ‘Korean AirBnB Script’. Moreover, the native Korean speakers who tried the website said that model struggled slightly with foreign names and domain specific words. From the feedback we got the model may benefit from more training data for non-korean words written with Hangul characters.

Model Variant	WER	CER	BLEU
Latest Model	<b>0.013</b>	<b>0.393</b>	<b>0.82</b>
LSTM	-	0.23	0.17

Table 2: The huge performance increase we got fine tuning transformers

#### 4.5 Error Analysis

Based on common errors made by our model we found out that:

1. The model struggles with severe obfuscations.
2. Foreign noun based hangul are error prone.
3. Domain-specific terms and proper nouns are more challenging to deobfuscate correctly.

### 5 Conclusions and Future Work

Our research demonstrates that a phonetic-aware sequence-to-sequence approach significantly improves text deobfuscation performance compared to traditional methods in the case of Korean "airbnb script". The use of phonetic representations helps bridge the gap between visually distinct but phonetically similar text, allowing our model to more effectively recover original text from obfuscated versions. Key findings from our work include:

- Phonetic preprocessing provides substantial benefits for text deobfuscation tasks.
- Character-level modeling is more effective than word-level approaches for this problem.
- Modern transformer architectures can effectively learn the complex mappings between obfuscated and original text forms.

Future work directions include:

1. Train on a even bigger dataset. We believe that this will be the most significant way to improve the performance of the model to be production ready.
2. Extending the approach to handle multiple languages and multilingual obfuscation. Potentially other language that uses phonetic alphabets, to incorporate G2P approach. This may also lead to a to foreign noun error we found with our models.
3. Incorporating visual similarity features to better handle cases where characters are substituted based on visual rather than phonetic similarity.
4. Developing adaptive models that can handle novel obfuscation techniques not seen during training.
5. Exploring applications in content moderation systems and accessibility tools. Potentially also assessing our model’s impact if any in the Korean online discourse after refinement and web deployment.

Our work contributes to the growing body of research on text normalization, for a very niche case of Korean text deobfuscation and provides practical tools for applications in content moderation, security, and accessibility.

### 6 Replicability

All the models and data used in our experiments are available at our GitHub repository: <https://github.com/Abdul03Rafay/KoreanDeobfuscator>.

The repository includes:

- Preprocessing scripts for phonetic conversion



- Model implementation using HuggingFace Transformers
- Training scripts
- Sample data (full data was too large to be uploaded.)
- Front-end functionality.

For full reproducibility, we also provide trained model checkpoints and detailed documentation on the environment setup.

## 7 Web Application

We've developed a website that bridges the communication gap for non-Korean speakers encountering obfuscated Korean text. Our intuitive three-step platform first allows users to experiment with Korean text obfuscation, giving them insight into how this process works. The core functionality lies in our transformer model that we have trained which deobfuscates 'Korean Airbnb Script'. For non-Korean users, we've integrated the Helsinki translation model, enabling translation of the deobfuscated text. This way we can help users not only decode obfuscated Korean text but also understand its meaning, making digital Korean content more accessible regardless of obfuscation techniques.



Figure 6: Bad Example With Partially Wrong Deobfuscation

## References

- [1] Limsopatham, N., & Collier, N. (2015). Adapting Phrase-based Machine Translation to Normalise Medical Terms in Social Media Messages. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pp. 1675-1680.
- [2] Bojja, N., Nedunchezian, A., & Wang, P. (2015). Machine translation in mobile games: augmenting social media text normalization with incentivized feedback. In *Proceedings of Machine Translation Summit XV: User Track*.
- [3] Lewis, M., Liu, Y., Goyal, N., Ghazvininejad, M., Mohamed, A., Levy, O., Stoyanov, V., & Zettlemoyer, L. (2020). BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 7871-7880.
- [4] Kumarage, T., Sheth, P., Moraffah, R., Garland, J., & Liu, H. (2023). How Reliable Are AI-Generated-Text Detectors? An Assessment Framework Using Evasive Soft Prompts. In *Findings of the Association for Computational Linguistics: EMNLP 2023*.
- [5] Honnet, P. E., Popescu-Belis, A., Musat, C., & Baeriswyl, M. (2018). Machine Translation of Low-Resource Spoken Dialects: Strategies for Normalizing Swiss German. In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*.
- [6] Tiwari, A. S., & Naskar, S. K. (2017). Normalization of Social Media Text using Deep Neural Networks. In *Proceedings of the 14th International Conference on Natural Language Processing (ICON-2017)*.



Figure 7: Good Example With Perfect Deobfuscation.