

**VIETNAM NATIONAL UNIVERSITY HO CHI MINH CITY
HO CHI MINH UNIVERSITY OF TECHNOLOGY
FACULTY OF COMPUTER SCIENCE AND ENGINEERING**



GRADUATION THESIS

**CONTEXTUAL VIETNAMESE SPELLING
CORRECTION**

Major: COMPUTER SCIENCE

THESIS COMMITTEE: COMPUTER SCIENCE 1
SUPERVISOR: ASSOC PROF. QUẢN THÀNH THƠ
REVIEWER: PROF. NGUYỄN HỨA PHÙNG

—o0o—

STUDENT 1: ĐỖ ANH KHOA 1852471
STUDENT 2: NGUYỄN TRẦN HIỂU 1852370

HO CHI MINH CITY, 05/2022

Declaration Of Authenticity

We declare that this research is our own work, conducted under the supervision and guidance of Assoc. Prof. Quan Thanh Tho. The result of our research is legitimate and has not been published in any form prior to this. All materials used within this research are collected by various sources and are appropriately listed in the references section.

In addition, within this research, we also used the results of several other authors and organizations. They have all been aptly referenced.

In any case of plagiarism, I stand by my actions and will be responsible for it. Ho Chi Minh City University of Technology, therefore, is not responsible for any copyright infringements conducted within our research.

Ho Chi Minh City, May, 2022

Author

Do Anh Khoa

Nguyen Tran Hieu

Acknowledgment

We are using this opportunity to express our gratitude to everyone who supported us during the making of our thesis. We are thankful for the invaluable knowledge, experience sharing received that had helped us improve.

We would like to show our gratitude to our mentor and supervisor throughout our thesis, Associate Professor Ph.D. Quan Thanh Tho. Throughout the making of this thesis, he has guided us with full support and constructive criticism.

We thank the Ho Chi Minh University of Technology for providing us the many great invaluable knowledge and lessons within four years of attending the school. Without attending the school, we would not have the chance to be friends and this thesis would not be possible.

Finally, our motivation to pursue this thesis comes from the support from our families and friends. We would like to show our deepest and most sincere gratitude to our parents, who sacrifice and support, and were with us in our life journey; our friends and siblings, whose encouragements and advises we took to heart in along the way.

Abstract

Spelling errors are common in everyday life, from writing emails to editing documents. The task consists of two tasks, namely detection to identify the error, and correction to suggest the right replacement. There are two main types of spelling errors: non-word and real-word errors. Non-word error is when the error itself is not in the dictionary and is not a known word, for example, "tôi" mistyped to "tooi". Real-word error is due to misspelling a word to make another word that is in the dictionary; for example, "tôi" mistyped to "toi", a diacritic error, is a real-word error. Real word error are harder to detect, as we have to rely on the surrounding context to perform reasoning and deduce the correct token. Moreover, spelling error are increasingly difficult when multiple tokens are merged together, for instance, "học hành" when removing space, becomes "họchanh", or when a token is split to sub-tokens like "tôi" to "t ôi".

Historically, rule-based methods are used to correct misspellings, which does not capture real-word errors and handle more complex misspelling cases, leading to studies of statistical models like N-gram. These model are trained on large corpus, and can capture information using the surrounding text. Nonetheless, a major setback is when two (or more) spelling mistakes stand near each other, the model will not be able to understand and capture the context correctly, which has been studied in [1]. Moreover, these method does not tackle the merged tokens and split token scenarios. Recently, two research group from Vietnam National University and Zalo Corporation propose contextual spelling correction using a deep-learning based approach.

In this paper, we study the effect of leveraging a pre-trained Transformer Encoder such as PhoBERT for Contextual Vietnamese Spelling Correction. We also adapt a tokenization repair module for Vietnamese, to handle merged syllable and split tokens scenario. With both model combined, we achieved mixed results on two public testing set by Vietnamese National University and Zalo Corporation.

Contents

1	Introduction	1
1.1	Motivation	1
1.2	Problem Statement	1
1.3	Goals	2
1.4	Scope	2
1.5	Thesis Structure	2
2	Theoretical Background	3
2.1	N-Grams	3
2.2	Minimum Edit distance	4
2.3	Neural networks	4
2.3.1	Feed-forward neural networks	5
2.3.2	Recurrent neural networks	6
2.3.3	Sequence-To-Sequence Architecture	10
2.3.4	Tokenization	12
2.4	Transformer Architecture	13
2.4.1	Attention	13
2.4.2	Self-attention	15
2.4.3	Multi-head Attention	17
2.4.4	Model Description	19
2.5	BERT Model	20
2.5.1	Definition	21
2.5.2	Model Architecture	21
2.5.3	Pre-Training BERT	23
2.5.4	Fine-tuning BERT	23
2.5.5	BERT Variants	24
3	Related Works	25
3.1	Transformer-based Model for Vietnamese Spelling Correction	25
3.1.1	Model	25
3.1.2	Testing set	26
3.1.3	Result and Evaluation	26
3.2	Hierarchical Transformer Encoders for Vietnamese Spelling Correction	27
3.2.1	Model	27
3.2.2	Testing set	28
3.2.3	Result and Evaluation	28
3.3	Soft-masked BERT	28
3.3.1	Detector	29

3.3.2	Corrector	29
4	Model	31
4.1	Model Architecture	31
4.1.1	BERT embedding	32
4.1.2	Detector	32
4.1.3	Corrector	32
4.1.4	Objective function	32
4.2	Tokenization Repair	33
4.2.1	Unidirectional Language Models	33
4.2.2	Bidirectional Sequence Labeling Model	34
4.3	Pipeline	35
5	Application Implementation	37
5.1	Programming Languages	37
5.1.1	Python	37
5.1.2	JavaScript	38
5.2	Tools	39
5.2.1	Google Colaboratory	39
6	Result and Evaluation	40
6.1	Dataset	40
6.1.1	Preprocessing Data	40
6.1.2	Errors and Augmenters	40
6.1.3	Augmentation Pipeline	41
6.1.4	Overall Data	42
6.2	Real world data analysis	43
6.3	Baselines and Model	44
6.3.1	Baselines models	44
6.3.1.1	Transformer	44
6.3.1.2	Soft-masked BERT	45
6.3.1.3	Hard-masked XLMR	45
6.3.2	Configuration settings	46
6.3.2.1	Transformers	46
6.3.2.2	Soft-masked BERT and Hard-masked XLMR	46
6.3.2.3	BERT Fine-tuned	47
6.3.3	Overall training settings	47
6.4	Evaluation metrics	47
6.4.1	Detection and Correction	47
6.4.2	Tokenization Repair	48
6.5	Result And Evaluation	48
6.6	Effect of Hyperparameter	52
7	Discussion and Future Work	53
8	Summary	57

List of Figures

2.1	A simple neural network	5
2.2	Single-layer Perceptron vs Multi-layer Perceptron	6
2.3	LSTM architecture	7
2.4	LSTM forget gate	8
2.5	LSTM input gate	8
2.6	LSTM cell state	9
2.7	LSTM cell state	9
2.8	Bidirectional Recurrent Neural Network	10
2.9	An overview of Sequence to Sequence architecture	11
2.10	Sequence to sequence model with Bahdanau attention mechanism.	14
2.11	Self-attention: Each word is associated with other words in different positions	15
2.12	Scale dot product attention	18
2.13	Multi-head self-attention	19
2.14	Transformer architecture	19
2.15	A sample of BERT architecture with multiple layer of Transformer block	21
2.16	Differences between BERT, OpenAI GPT and ELMo architecture	21
2.17	BERT's input representation	22
2.18	Overall pre-training and fine-tuning procedures for BERT.	23
2.19	The many spawned variants of BERT	24
3.1	VSEC model	25
3.2	Hierarchical Transformer Model	27
3.3	Soft-masked BERT.	28
4.1	Our proposed BERT fine-tuned model	31
4.2	Unidirectional Language Model with LSTM	33
4.3	Bidirectional Sequence Labelling with BiLSTM	34
4.4	Spelling correction pipeline	35
5.1	Web-app interface to demonstrate result.	38
5.2	Choose a model to demonstrate performance	39
6.1	Data augmentation pipeline	42
6.2	Errors type percentages	43
6.3	Percentages of error tokens	44
6.4	Hard-masked BERT.	45

List of Tables

4.1	Example of pipeline	36
6.1	Example of some misspelling errors	41
6.2	Data for spelling correction	43
6.3	Real world data counts	43
6.4	Results	49
6.5	Example of correction	50
6.6	VSEC test set	50
6.7	Viwiki test set	51
6.8	Tokenization repair result	51
6.9	Example of tokenization repair	52
6.10	Example of correction with tokenization repair	52
6.11	Effect of training set size on accuracy	52
7.1	Example of domain-specific keyword	54
7.2	Example of English word and Shortcut	55

Chapter 1

Introduction

1.1 Motivation

We as human will inevitably make mistake during typing, whether it is accidental key-strokes or misunderstanding of its meaning. As we need our emails and document to be error-free and presentable, spelling correction is vital and has a wide application value to our life. Although this problem is extensively researched in English [2, 3, 4] and Chinese [5, 6], it is still relatively new for Vietnamese.

Previous studies approach the problem using a rule-based methods, which works for simplex structured cased. However, it often fails to correct harder misspelling cases like close proximity keystroke and is unable to detect real-word cases where word have no visible error due to the method not utilizing context. Statistical model approach like n-gram[7, 8] yield a more promising result, able to capture the surrounding information to deduce the correct answer. Nonetheless, it is fragile to consecutive spelling errors, which can break the context and reduce accuracy. In November 2021, Zalo research team release a paper on spelling correction task by proposing a hierarchical transformer encoder[1], which combine both character and word level representation, show promising results. At the same time, Vietnamese National University release a different approach based on machine translation aproach, translating an erroneous sentence to an error free sentence. [9]

Spelling correction has not been widely studied in Vietnamese compared to Chinese and English, although its application is important. In real applications, search engines can utilize spelling correction systems to improve user experience. Many word processing software is working on building spelling correction systems for day to day usages such as Microsoft Word, Google Docs, or third party integration like Grammarly for English. It can also be generalised to error correction task, which can be useful in other field such as Optical Character Recognition, which studies the textual extraction of image, or Automatic Speech Recognition,

1.2 Problem Statement

For an input text sequence $X = [x_1, x_2, \dots, x_n]$, each token may or may not contains spelling error, our goal is to generate the correct spelling of the corresponding text $Y = [y_1, y_2, \dots, y_n]$. Each token y_i at position i is either retained from x_i , replaced as another token, or a special $\langle del \rangle$ token for removal. Note that the input and output sequence are of the same length, however after

apply special `` token, the final result might have different number of tokens.

1.3 Goals

In this thesis, we aim to achieve several outcomes:

- Performing survey on Vietnamese spelling errors.
- Study the effectiveness leveraging a pre-trained Transformer Encoder such as phoBERT for the Contextual Vietnamese Spelling Correction task.
- Enhance our correction model with a tokenization repair module to handle merged syllables and split tokens cases.
- Perform a comparative study between different deep-learning approaches.
- Evaluate our approach against published testing set from other researches.
- Build a simple web application to demonstrate our results.

1.4 Scope

Per the scope of our thesis, we will focus only on Vietnamese spelling errors. This exclude foreign language correction, shortcuts and teencodes which we believed to be an intentional input.

1.5 Thesis Structure

There are totally eight chapter in this thesis proposal:

Chapter	Content
1	A brief introduction about plan and objectives of thesis
2	Introduction of theoretical background as foundation knowledge that are applied in the project
3	Related works have been done on this problem
4	Solution and design approach for problem statement of project
5	An overview of technology usage and application
6	Results and test set evaluation
7	Discussion and future works
8	Thesis summary

Chapter 2

Theoretical Background

2.1 N-Grams

An n-gram is a sequence of n words or tokens: a 2-gram (bi-gram) is a two-word sequence and a 3-gram (tri-gram) is a three-word sequence and so forth. It is a probabilistic model that could assign probability of the last word of n-gram given the previous words, and also to assign probabilities to the whole sequence, which is called **language model**.

To compute the probability of an entire sequence with n words $P(w_1, w_2, \dots, w_m)$, traditionally, chain rule is applied:

$$\begin{aligned} P(w_1, w_2, \dots, w_m) &= P(w_1)P(w_2|w_1)P(w_3|w_1w_2)\dots P(w_m|w_1\dots w_{m-1}) \\ &= \prod_{i=1}^m P(w_i|w_1\dots w_{i-1}) \end{aligned}$$

The intuition of the n-gram model is that the probability of a word is calculated as a conditional probability given its last n words rather than its whole history. For instance, a bi-gram model could approximate the probability of a word given its previous word

$$P(w_m|w_1\dots w_m) \approx P(w_m|w_{m-1})$$

This approximate assumption is called a **Markov** assumption. From that assumption, we could generalize n-gram by looking at $n - 1$ previous word. Thus, the general equation for n-gram approximation to the conditional probability of the next word in a sequence is:

$$P(w_m|w_1\dots w_{m-1}) \approx P(w_m|w_{m-n+1}\dots w_{m-1})$$

Hence, the probability of the entire word sequence is for bi-gram is:

$$P(w_1\dots w_m) \approx \prod_{i=1}^m P(w_i|w_{i-1})$$

Maximum likelihood estimation (MLE) is used for estimating these probabilities. We get an MLE estimate for parameters of the n-gram model by getting counts from a corpus and then

normalizing those counts to lie between 0 and 1.

$$P(w_m | w_{m-n+1} \dots w_{n-1}) = \frac{C(w_{m-n+1} \dots w_{m-1} w_n)}{C(w_{m-n+1} \dots w_m)}$$

The estimation of n-gram probability is computed by dividing the observed frequency of a particular sequence by the observed frequency of a prefix, this ratio is called relative frequency.

2.2 Minimum Edit distance

In computational linguistics and computer science, edit distance is a way of quantifying how dissimilar two strings (e.g., words) are to one another by counting the minimum number of operations required to transform one string into the other. There are several algorithms that have been used to solve this problem, though Levenshtein Distance (V. Levenshtein, 1965) is more widely used in the field for error correction.

Levenshtein distance is a number that tells you how different two strings are, the higher the number, the more different the two strings are. The distance is measured by a set of three operations on single-character edits: insertion, deletion and substitution. For instance, the string "kitten" and "sitting" is 3 edit distance away due to:

- "kitten" → "sitting": substitution of "s" for "k"
- "sitten" → "sittin": substitution of "e" for "i"
- "sittin" → "sitting": insertion of "g"

Given two string **a**, **b** of length $|a|$ and $|b|$, then $\text{lev}(\mathbf{a}, \mathbf{b})$ is calculated by the following function:

$$\text{lev}_{a,b}(i, j) = \begin{cases} \max(i, j) & \text{if } \min(i, j) = 0 \\ \min \begin{cases} \text{lev}_{a,b}(i-1, j) + 1 \\ \text{lev}_{a,b}(i, j-1) + 1 \\ \text{lev}_{a,b}(i-1, j-1) + 1_{(a_i \neq b_j)} \end{cases} & \text{otherwise} \end{cases}$$

where i, j are character position of string **a** and **b**, $1_{(a_i \neq b_j)}$ denotes 0 when $a = b$ and 1 otherwise.

Although using minimum edit distances would eliminate several options for choosing the appropriate word, there could still be a possibility where there are multiple words with the same edit distance values. Therefore, using the Levenshtein similarity ratio could be useful in such cases:

$$\text{ratio} = \frac{(|a| + |b|) - \text{lev}_{a,b}(i, j)}{|a| + |b|}$$

where $|a|$ and $|b|$ are the length of sequence **a** and **b**, respectively.

2.3 Neural networks

Neural network is the basic concept of the field of Deep learning in general. Originally inspired by biological neural networks, this mathematical model mimics the interconnected network of

neurons transmitting elaborate patterns of signal. The image below illustrates a simple neural network.

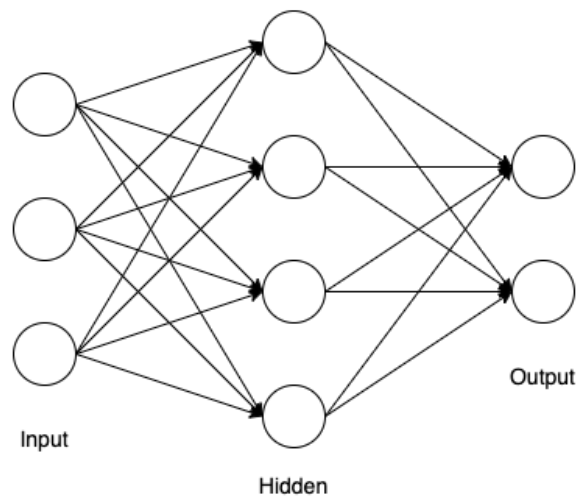


Figure 2.1: A simple neural network

A neural network is composed of a large number of different types of neurons. The image illustrates three types: input layers consist of 3 neurons, hidden layers consist of 4 neurons, and output layer consists of 2 neurons. The hidden neurons, or nodes, are where the computation is done and transfer the signal to the output layer. Neurons transmit their signal to others through a link called connection. In a sense, a neuron i "transmits" its signal to another neuron j through the connection, and each connection holds a weight w_{ij} value that decides how much of the signal is transferred.

In the deep learning field, neural networks have been applied to various tasks. Natural language processing tasks such as text classification, tokens classification (part-of-speech tagging, name-entity recognition),... witness a boost in performance by training different neural architecture. Training in deep learning in a sense is adjusting the weight and bias values to produce the desired output.

2.3.1 Feed-forward neural networks

A feed-forward network is an artificial network where the connections within the network do not form a cycle. Therefore, this type of network allows information to move in only one forward direction, from the input node and ends at output nodes.

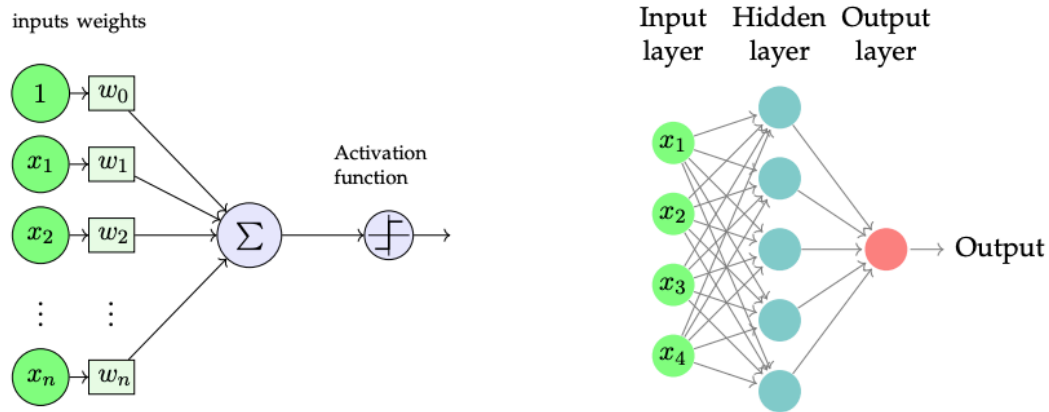


Figure 2.2: Single-layer Perceptron vs Multi-layer Perceptron

Single-layer Perceptron

Single layer network does not contain any hidden layers. It consists of an input layer and output layer of nodes. In neural nets, we do not account for the input layer as it only receive input fed from the outside, while no computation is done. Let us assume $\mathbf{x} = (x_1, x_2, \dots, x_n)$ as a vector of dimension n to be fed into the network. Then, output z , the sum weight of inputs have the formula:

$$z = \sum_{i=1}^n w_i x_i + b = \mathbf{w}^T \mathbf{x} + \mathbf{b}$$

where $\mathbf{w} = (w_1, w_2, \dots, w_n)$ is a vector containing weight values, b is the bias value. This linear transformation z is then passed through an activation function to produce output.

Multi-layer perceptron

An extension to single-layer perceptron, this class of network contains more than 1 layer for more complex computation, still in a direct way. Each node in a layer has at least 1 connection to a node in the subsequent layer. These networks may contain more than one hidden layer, which allows them to learn the non-linear representation of the data presented.

2.3.2 Recurrent neural networks

Sequential modeling is the study of generating a sequence of values by analyzing a series of input values. The ability to handle and process sequences are important in the field of Natural Language Processing; for example, given a text (sequence of words or characters), we want to analyze the sentiment values associated with it. With the traditional feed-forward model, we cannot represent data with various timestamps to learn effectively.

Definition

A recurrent neural network (RNN) is a type of neural network that can handle sequential data effectively. An extension of traditional neural network possessed by RNN is the cyclic connection that each unit can have with each other; while simple ANN network learns to produce an output from a given input, RNN can utilize the inputs and outputs of all previous units, thus providing memory-like internal states for the network.

Let us define a sequence of inputs $\mathbf{X} = (x_1, x_2, \dots, x_n)$ that is compressed into a fixed dimension vector. This allows us to deal with variable-length input and outputs. At timestep t , we have a vector h_{t-1} that contains the hidden state value up to $t - 1$ timestep. The recurrent neural network will compute the new internal state h_t from the previous internal state and current input as follow:

$$h_t = \text{activation}(W_h x_t + U_h h_{t-1} + b_h)$$

$$\hat{y}_t = \text{softmax}(W_y h_t + b_y)$$

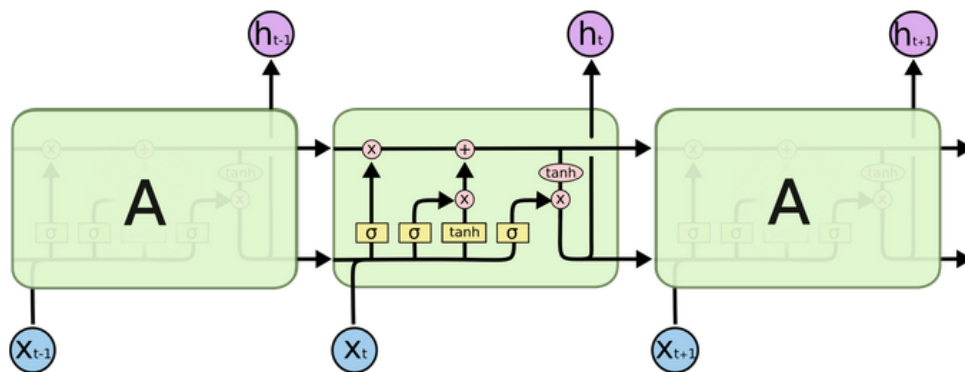
where W_h is the weight matrix of input, U_h is the recurrent weight matrix, and b_h is the bias values. After we obtain the hidden value for the current timestep, we calculate the output y_t . Training a recurrent neural network refers to adjusting the parameters in these equations.

In RNNs, we encounter a phenomenon called *vanishing gradient*. The multiplicative gradient can be exponentially increasing or decreasing, hence RNNs will suffer from capturing long-term dependencies. This is a major problem as long-term dependencies are one of the main characteristics of language.

Efforts have been made to confront the vanishing gradient problem. Most notably, to architecture called Gated Recurrent Unit (GRU) [10] or Long Short-term memory (LSTM) [11] architecture. Within the scope of the thesis, we chose the latter to implement, which we will discuss in the next section.

Long short-term memory

Long Short-Term Memory (LSTM) [11] network is an architecture of Recurrent Neural Network, as it is used to resolve long-term dependencies as stated above. Today, it is one of the most widely used models in Deep Learning for NLP.



The repeating module in an LSTM contains four interacting layers.

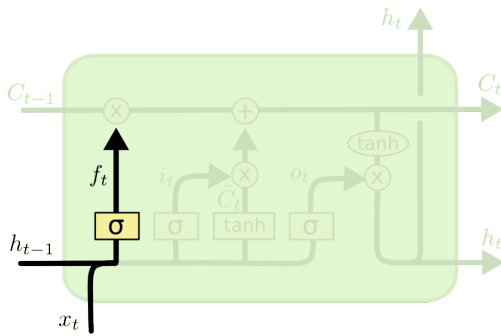
Figure 2.3: LSTM architecture

Source: <https://colah.github.io/posts/2015-08-Understanding-LSTMs/>

The heart of LSTM network is the cell state and gate units that control the internal states. The gate unites consist of three units: input gate, forget gate, output gate.

Forget gate

The forget gate mechanism allows LSTM to decide the information that will be kept or discarded when passed through this timestep. It uses a sigmoid function layer called the "forget gate layer". The input is a combination of current timestep x_t and previous hidden state h_{t-1} , and the sigmoid function returns a value ranging from 0 to 1, which resembles how much of the information will be retained. In this case, 1 stands for completely retaining information, while 0 stands for the opposite.



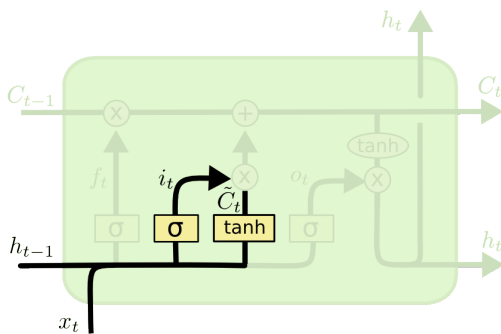
$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$

Figure 2.4: LSTM forget gate

Source: <https://colah.github.io/posts/2015-08-Understanding-LSTMs/>

Input gate

To make the decision on what information will be stored in the cell state, the input gate passes input through two different activation functions. Firstly, a sigmoid layer decide which value will be updated, which also outputs a value from 0 to 1. At the same time, the tanh layer creates a vector of new candidate values to be added to the state.



$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

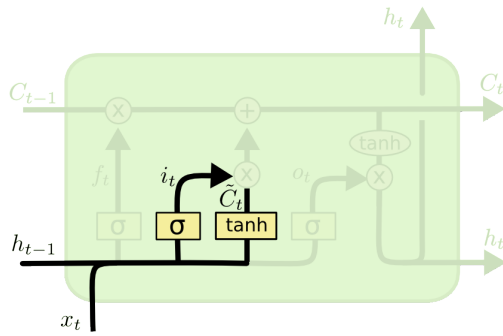
$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

Figure 2.5: LSTM input gate

Source: <https://colah.github.io/posts/2015-08-Understanding-LSTMs/>

Cell state

To transform the old cell state C_{t-1} to the new cell state, we firstly perform "forgetting" the input we decided at the forget gate step, through a point-wise multiplication operation. Then, we add on the new information from the input gate to achieve the current cell state through a point-wise addition operation.



$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

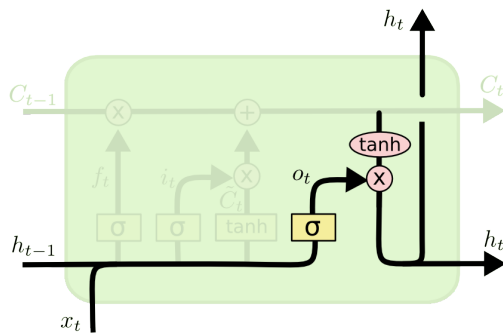
$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

Figure 2.6: LSTM cell state

Source: <https://colah.github.io/posts/2015-08-Understanding-LSTMs/>

Output gate

Finally, we achieve the output through the output gate, which takes in the previous hidden state and current timestep input, and the cell state. Through the sigmoid layer, LSTM decides the part of the cell state to be the output. Then combine it with the cell state that has been through tanh and produce output through point-wise multiplication.



$$o_t = \sigma(W_o [h_{t-1}, x_t] + b_o)$$

$$h_t = o_t * \tanh(C_t)$$

Figure 2.7: LSTM cell state

Source: <https://colah.github.io/posts/2015-08-Understanding-LSTMs/>

Bidirectional Recurrent Neural network

The previous discussion focus on RNNs that are able to observe past input and hidden states in order to learn and predict the next symbol in a data sequence. However, for many sequence labeling tasks, we would also like to have access to future context. To achieve this, we can use two RNN models that read the input sequence in two ways: right-to-left propagation and left-to-right propagation. This concept is commonly referred to as Bidirectional Recurrent Neural networks (BRNNs) [12].

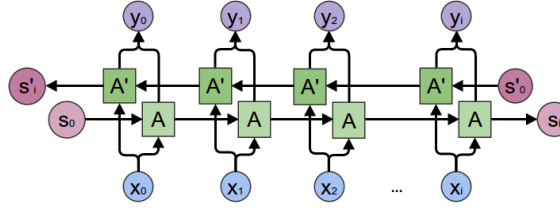


Figure 2.8: Bidirectional Recurrent Neural Network

For Bi-RNN, the right-to-left propagation is the same as the traditional RNN, while the left-to-right propagation runs in the opposite direction. As a result, the input sequence is forwarded to 2 different hidden layers for computation.

$$\vec{h}_t = \tanh(\vec{W}x_t + \vec{U}\vec{h}_{t-1}) + \vec{b}$$

$$\overleftarrow{h}_t = \tanh(\overleftarrow{W}x_t + \overleftarrow{U}\overleftarrow{h}_{t+1}) + \overleftarrow{b}$$

Note how the calculation of right-to-left propagation hidden state depends on the hidden state of future timestep. We will concatenate the hidden state of both right-to-left \overleftarrow{h}_t and left-to-right \vec{h}_t to obtain the hidden state h_t that will forward to the output layer.

Bidirectional recurrent neural networks have improved results in various domains [13, 14].

2.3.3 Sequence-To-Sequence Architecture

Sequence-to-Sequence (Sutskever et al.2014) [15] is a Deep Neural Networks model with the goal of generating output sequence from an input sequence whose length may be different from that of output. Although the model was initially proposed to tackle Machine Translation problem, Sequence-to-Sequence has also been applied in numerous other systems such as Speech recognition, Text summarization, etc.

The model comprises 2 main parts: Encoder and Decoder, both are recurrent neural networks that are connected through a context vector.

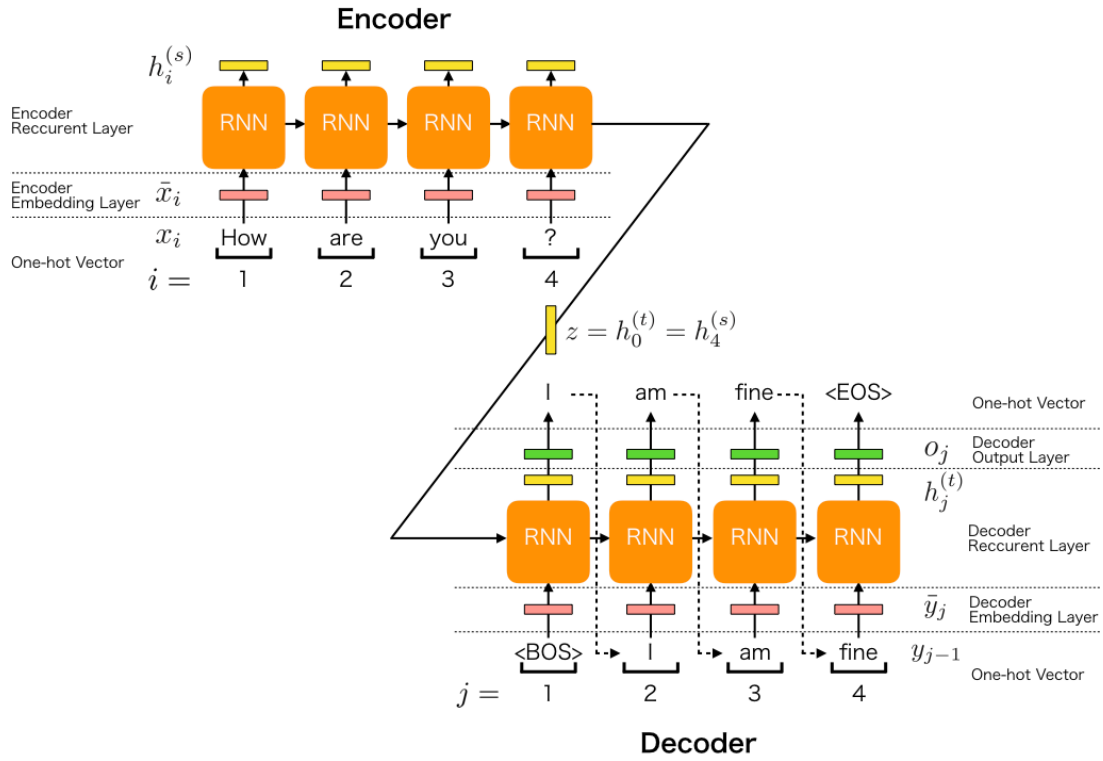


Figure 2.9: An overview of Sequence to Sequence architecture

Image source: Tokui, Seiya, et al. [16]

Encoder

This encoder is a stack of recurrent units (which could be LSTM (Long short-term memory) or GRU (Gated recurrent unit) to avoid vanishing gradient). It is fed in overtime an input sequence $\mathbf{x} = x_1, x_2, \dots, x_n$ with x_t is the t^{th} token of the sequence, processes each token of the input and tries to encapsulate all the information from the input sentence and store them in its final internal state, which is hidden state h_t calculated by the following formula:

$$h_t = f(x_t, h_{t-1}) = f(W^{(hh)}h_{t-1} + W^{(hx)}x_t)$$

These final hidden states, called *context vector*, are then passed onto the decoder part and act as an initial hidden state for the decoder block to produce the target sentence.

In order to boost the performance, unidirectional RNNs are replaced with the bidirectional due to their capabilities of capturing both forward and backward direction of the sequence.

Decoder

While the encoder packs all the information and encoded the sequence into a context vector, the decoder is used to decode that vector and generate the output sequence. The same structure is applied to the decoder block. A stack of several recurrent units that produce an output y_t for each time step t . Each recurrent unit accepts a hidden state from the previous unit, also produces and outputs its own hidden state. These hidden states are calculated by:

$$h_t = f(h_t) = f(W^{(hh)}h_{t-1})$$

A fully-connected layer is added as the output layer with the purpose of finding the probability vector to determine the final output along with the time steps by softmax.

$$y_t = \text{softmax}(W^s h_t)$$

Finally, the loss is calculated on the predicted outputs and then the errors are backpropagated to update the parameters. And through training over a period of time with sufficiently large data, the model could provide predictions with good performance.

2.3.4 Tokenization

Tokenization is a process of dividing a sequence into smaller units, called tokens, which then can be fed into a model. By tokenizing an unstructured stream of text, we achieve meaningful representations that are more useful for our model to learn. Choosing the right tokenization method is important to the spelling correction task. We present some tokenization methods commonly used in natural language processing tasks.

- **Word level tokenization:** In this method, sentences are split up into list of tokens separated by white-space. Each token will be monosyllabic in Vietnamese per the nature of the language. This technique is commonly used in word embedding methods like GloVe, word2vec.
- **Multi-word level tokenization:** Language like Vietnamese uses a composite number of monosyllabic tokens to convey an actual word. For example, "xác định" is composed of 2 syllables "xác" and "định", which have no relation to the composite word in meaning. Multi-word levels are useful to convey the meaning of the word, hence used in NLP tasks such as part-of-speech tagging, name-entity recognition. Some open-source libraries provide support for multi-word tokenization such as underthesea, pyvivn,...
- **Character level tokenization:** Word and above usually suffer from a large vocabulary size, which is especially true for spelling correction tasks where one syllable can have multiple different types of errors. Two problems arise: a large vocabulary size causes model to train and converge slower, and out-of-vocabulary problem for word that the model has never seen. Character tokens are atomic as characters are the basic unit that composes a word, hence a smaller set of tokens in vocabulary can fully cover a wide range of words, reducing out-of-vocabulary tokens. A word embedding method known as fasttext utilize character tokenization.
- **Subword tokenization:** Although character tokenization addresses the problem of word-level tokenization, it is not without limitation. Dividing a sentence to character level greatly increase sequence length, which leads to slow training and inference time. Another problem is character level method cannot convey word context as well as word-level. Subword tokenization inherits both word and character level embedding. The idea is to observe the common patterns of multiple characters from words to form the subword units. For example, "faster" may be decomposed to "fast", "er".

Subword tokenization is commonly referred to as inbetween word and character level, and is used in many modern models in transformer-based architecture like BERT [17], or its variations. A popular algorithm used to train subword tokenizer is Byte pair encoding [18, 19]. The byte pair encoding algorithm works as follows:

- Step 1: Sentences are split into word level. Each word is appended a special token $\backslash w$ to mark the end of word.
- Step 2: Split all words into their character components and add to vocabulary. At this initial step, vocab only contains character tokens.
- Step 3: Count all occurrence of token pair that appears in the corpus, and add frequent pairs to the vocabulary as a new token.
- Step 4: Repeat step 3 until desired vocab size is reached or several merge operation is performed.

Training BPE is unsupervised and the only hyperparameter we can tune is vocab size or number of merge operations.

2.4 Transformer Architecture

The paper "Attention is all you need" [20], one of the most impactful paper in 2017, introduce the transformer architecture, which fully relies on attention mechanism, instead of the recurrent mechanism of traditional method like Recurrent Neural Network, Long Short-Term Memory,... To understand the proposed model, we first discuss attention.

2.4.1 Attention

Sequence-to-sequence models work great for machine translation tasks. However, its disadvantage is based on the fact that the neural network will compress all information from the sentence to a fixed-length vector. As Cho et al. (2014b). [21] has pointed out that the performance of a basic sequence-to-sequence deteriorates proportionally to the sequence length. Bahdanau et al. [22] propose the attention mechanism, an extension of encoder-decoder architecture to address this problem.

In an abstract level, attention extension does not build a single context vector out of the encoder's last hidden state. Instead, what makes attention mechanism special is that it create shortcuts between the context vectors and the entire source input. After that, the decoder can soft-search for the set of positions in the source sentence for the most relevant information. This is done by assigning customisable weighted sum of context vectors

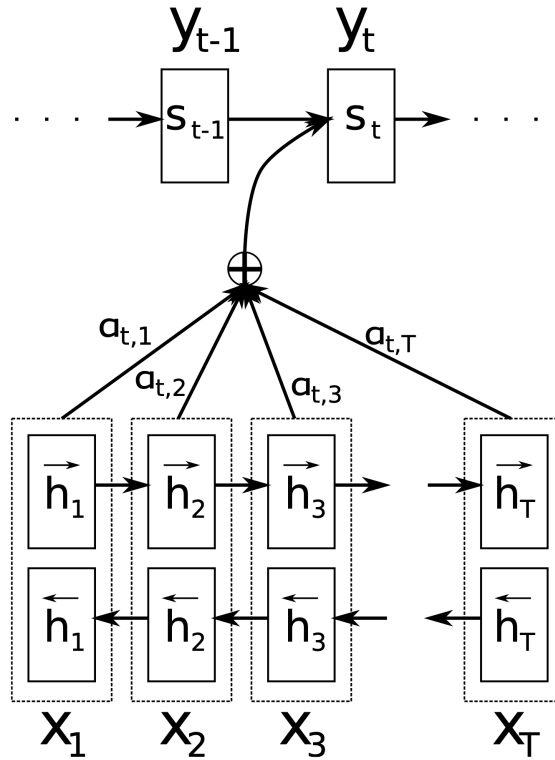


Figure 2.10: Sequence to sequence model with Bahdanau attention mechanism.

Image source: Bahdanau et al. [22]

Lets define the problem of machine translation with input source sequence $\mathbf{x} = (x_1, x_2, \dots, x_n)$ with length n , and output target sequence $\mathbf{y} = (y_1, y_2, \dots, y_m)$ with length m .

The encoder (which is a bidirectional RNN), contains the foward hidden states \vec{h}_i and backward hidden state \overleftarrow{h}_i . The overall representation is denoted as

$$\mathbf{h}_i = [\vec{h}_i^T, \overleftarrow{h}_i^T]^T, i = 1, 2, \dots, n$$

In the proposed model, decoder has hidden state $\mathbf{s}_i = f(\mathbf{s}_{i-1}, y_{i-1} \mathbf{c}_i)$. \mathbf{c}_i denotes for the context vector at position i , which is calculated by the weighted sum of hidden state of the input sequence \mathbf{x} . The weight is an alignment score, calculate as below:

$$\mathbf{c}_t = \sum_{j=1}^n \alpha_{t,j} \mathbf{h}_j$$

$$\alpha_{t,j} = \frac{\exp(\text{score}(s_{t-1}, h_j))}{\sum_{k=1}^n \exp(\text{score}(s_{t-1}, h_k))}$$

Here, the alignment model will score how matching the input at position j and output at position t match. In the original paper, the alignment model will be parameterized as a feedforward neural network which will be jointly learned with the model.

$$\text{score}(\mathbf{s}_t, \mathbf{h}_j) = \mathbf{v}_a^T \tanh(\mathbf{W}_a [\mathbf{s}_t; \mathbf{h}_j])$$

where \mathbf{v} and \mathbf{W} are learnable weight matrices. Since then, multiple studies have extended the

attention mechanism [23, 20]. Sequence-to-sequence models with attention improve the performance of traditional models in machine translation tasks and have achieved an increase in translation performance.

2.4.2 Self-attention

Self-attention, also known as intra-attention, is an attention mechanism relating different positions of a single sequence in order to compute a representation of the same sequence. For each word, this mechanism allow the model to find different positions in sequence for relevancy. Self-attention has been adopted by some papers [24]. The term "self-attention" was coined when the paper "Attention is all you need" [20] is published.

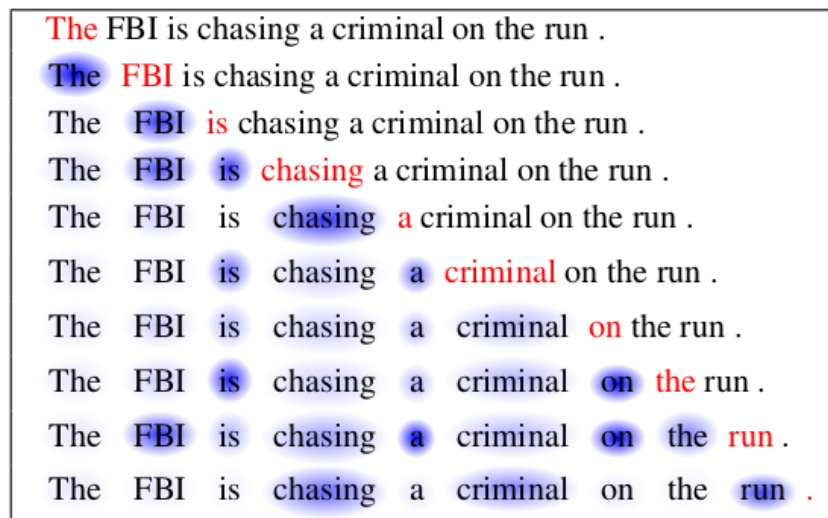


Figure 2.11: Self-attention: Each word is associated with other words in different positions

Image source: Cheng et al. [22]

The procedure for self attention is as follows:

- Step 1: Perform input embeddings.
- Step 2: Derive from weight matrices a set of **Key**, **Value**, **Query** for each embedding.
- Step 3: Calculate attention scores for each input by multiplying each **Query** vector to all **Key** vector
- Step 4: Calculate softmax of attention scores
- Step 5: Multiply attention scores with **Values**
- Step 6: Calculate the sum weighted for each output.

To better understand the mechanism, lets perform a simple example.

First, lets initialize 3 input embeddings with dimension 4, and weight matrices corresponding to Key, Value, Query. Assuming Key, Value, Query are all of dimension 3, then these weight matrix will have dimension of 4x3.

```
## Initialize input embedding
Input1 = [1,0,0,1]
Input2 = [0,2,2,0]
Input3 = [2,1,1,2]

## Initialize key weight matrix
Weight_key = [[0, 0, 1],
              [1, 1, 0],
              [0, 1, 0],
              [1, 1, 0]]

## Initialize query weight matrix
Weight_query = [[1, 0, 1],
                [1, 0, 0],
                [0, 0, 1],
                [0, 1, 1]]

## Initialize value weight matrix
Weight_value = [[0, 2, 0],
                [0, 3, 0],
                [1, 0, 3],
                [1, 1, 0]]
```

Step 2: Calculate, **key**, **value** and **query** representation for input embedding.

```
# Input 1's key, value, query calculation
### Key
      [[0, 0, 1],
[1,0,0,1] x [1, 1, 0], = [1,1,1]
      [0, 1, 0],
      [1, 1, 0]]

### Value
      [[0, 2, 0],
[1,0,0,1] x [0, 3, 0], = [1,3,0]
      [1, 0, 3],
      [1, 1, 0]]

### Query
      [[1, 0, 1],
[1,0,0,1] x [1, 0, 0], = [1,1,2]
      [0, 0, 1],
      [0, 1, 1]]

*** We can see that these calculations can operate with multiple embedding
    at the same time.
### Key
      [[0, 0, 1],
[[1,0,0,1], x [1, 1, 0], = [[1,1,1],
[0,2,2,0],   [0, 1, 0],   [2,4,0],
[2,1,1,2]]   [1, 1, 0]]   [3,4,2]]

### Value
      [[0, 2, 0],
[[1,0,0,1], x [0, 3, 0], = [[1,3,0],
[0,2,2,0],   [1, 0, 3],   [0,6,0],
[2,1,1,2]]   [1, 1, 0]]   [3,9,3]]

### Query
      [[1, 0, 1],
```



```
[[1,0,0,1], x [1, 0, 0], = [[1,1,2],  
[0,2,2,0],   [0, 0, 1],   [2,0,2],  
[2,1,1,2]]   [0, 1, 1]]   [3,2,5]]
```

Next, let's calculate attention scores for each input and apply softmax function

```
### Input1  
                T  
                [[1,1,1],  
[1,1,2] x [2,4,0], = [4, 6, 18]  
          [3,4,2]]  
  
softmax([4,6,18]) = [0.0, 0.1, 0.9]
```

Next, we multiply scores with values and get the sum weighted as attention representation for input 1

```
### multiply score with value  
0.0 x [1, 3, 0] = [0, 0, 0]  
0.1 x [0, 6, 0] = [0, 0.6, 0]  
0.9 x [3, 9, 3] = [2.7, 8.1, 2.7]  
  
### attention representation for Input1  
  [0 , 0 , 0 ]  
+ [0 , 0.6, 0 ]  
+ [2.7, 8.1, 2.7]  
=  [2.7, 8.7, 2.7]
```

We can repeat this process to achieve attention representation for the other 2 embeddings.

2.4.3 Multi-head Attention

What makes the transformer model special is its multi-head self-attention mechanism. We take a deep dive to see what the attention is all about.

Scaled dot-product attention

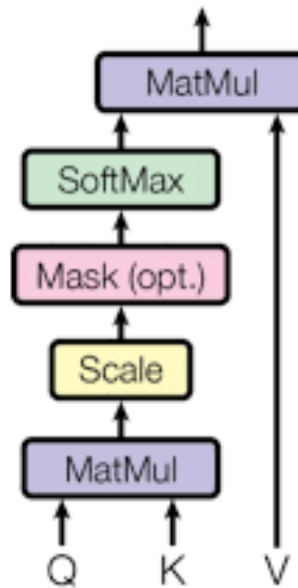


Figure 2.12: Scale dot product attention

Image source: Vaswani et al. [20]

The scaled dot-product attention mechanism adopts the general attention. Firstly, we calculate the dot product between query Q and key K . The result is then scaled by dividing by a factor of $\sqrt{d_k}$ and proceed to apply the softmax function. The result is the weight used to scale the value V .

$$attention(Q, K, V) = \frac{QK^T}{\sqrt{d_k}}V$$

Multi-head Attention

The multi-head mechanism allow us to run multiple scaled dot product attention heads in parallel. After the attention representation have been calculated, it is concatenate and sip through an affine layer to “allows the model to jointly attend to information from different representation subspaces at different positions. With a single attention head, averaging inhibits this.” [20].

$$MultiHead(Q, K, V) = Concat(head_1, \dots, head_h)W^O$$

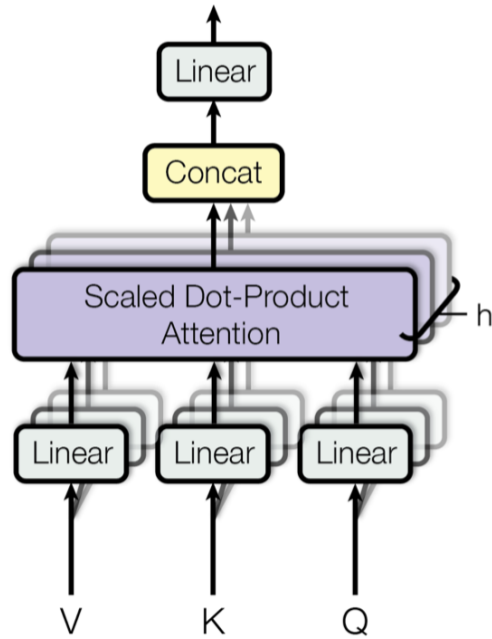


Figure 2.13: Multi-head self-attention
*Image source: Vaswani et al. [20]

2.4.4 Model Description

The architecture of transformer is as follow:

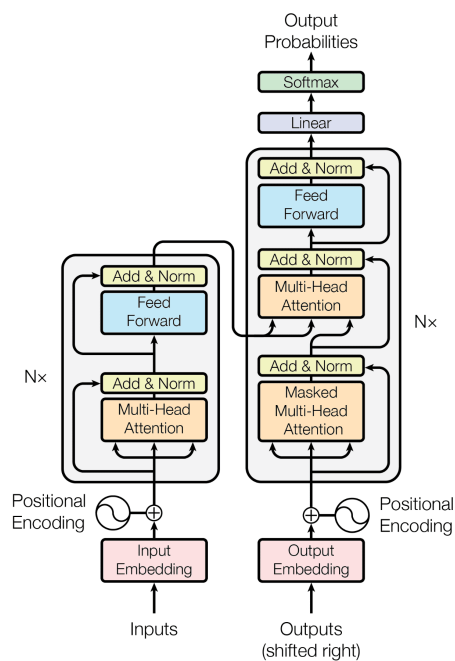


Figure 2.14: Transformer architecture

Encoder

Input embedding

The tokenized input will be fed into the input embedding layer to extract the vector representation.

Positional encoding

The input embedding encodes the context and meaning of a word. However, at different positions, words may convey a different meaning. This architecture adds a layer to encode the positional information relating to each word.

$$PE_{(pos, 2i)} = \sin(pos/10000^{2i/d_{model}})$$

$$PE_{(pos, 2i+1)} = \cos(pos/10000^{2i/d_{model}})$$

where pos is the relative position of each word, and PE is the positional encoding at position i . The sum of positional encoding and input embedding will be fed to the actual encoder layers.

Encoder layers

The encoder generates an attention-based representation with the capability to locate a specific piece of information from a potentially infinitely-large context. The original paper use 6 layers of encoding, each comprised of multi-head attention and a feed-forward layer. The first layer receives input + positional embedding, while others receive the output of the previous layers.

Decoder

Similar to the encoder, the decoder comprises N layers. The first layer receives the input and positional embedding of the target sentence, while the other layer receives output from the previous layer. An important difference is an introduction of masking multi-head attention. By introducing masking, the decoder only attends to input embedding precedes it, whereas the encoder can attend to both directions. The masking makes the decoder unidirectional.

Encoder-Decoder attention

In addition, another multi-head attention is put in the decoder. In these layers, the queries come from the previous decoder layer, and the memory keys and values come from the output of the encoder. This enables the decoder to attend to all positions of the source input.

2.5 BERT Model

There have been many state-of-the-art NPL proposals since the introduction of pre-trained language models such as ELMo (Peters et al., 2018a) [25], Generative Pre-trained Transformer (OpenAI GPT) (Radford et al., 2018) [26], ULMFit (Howard, Ruder et al., 2018) [27]. Downstream tasks could be solved using the two following approaches for applying pre-trained language representations:

- Feature-based: this approach address the usage of task-specific architectures that based on pre-trained representations as additional features (e.g., ELMo)
- Fine-tuning: task-specific parameters are introduced and trained on downstream tasks by fine-tuning all pre-trained parameters (e.g., OpenAI GPT, ULMFit)

These models refrain from reaching the actual power of the pre-trained model due to its unidirectional nature, which limits the option of architectures that could be used during pre-training. BERT is proposed to tackle this problem.

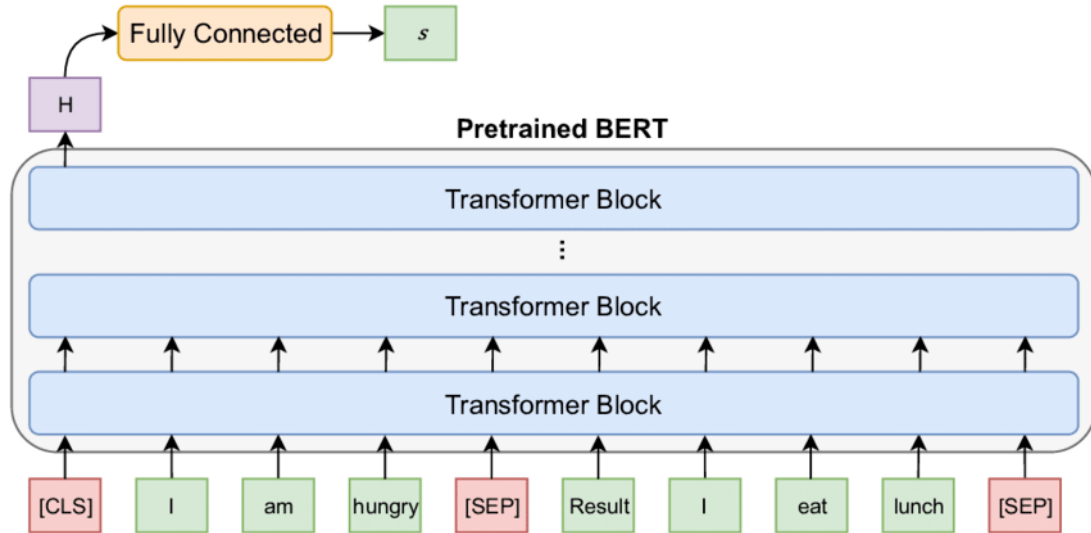


Figure 2.15: A sample of BERT architecture with multiple layer of Transformer block

2.5.1 Definition

BERT (Bidirectional Encoder Representations from Transformers) (Jacob Devlin et al., 2018) [17] is a transformer-based machine learning model developed with fine-tuning approach by researchers at Google AI Language. BERT is designed to overcome the limitations that persist in already pre-trained models such as OpenAI GPT and ELMo while also enhancing the ability to encapsulate the context of sentences. Another distinguishing feature of BERT is that there are minimal changes in the pre-trained architecture and final downstream architecture.

2.5.2 Model Architecture

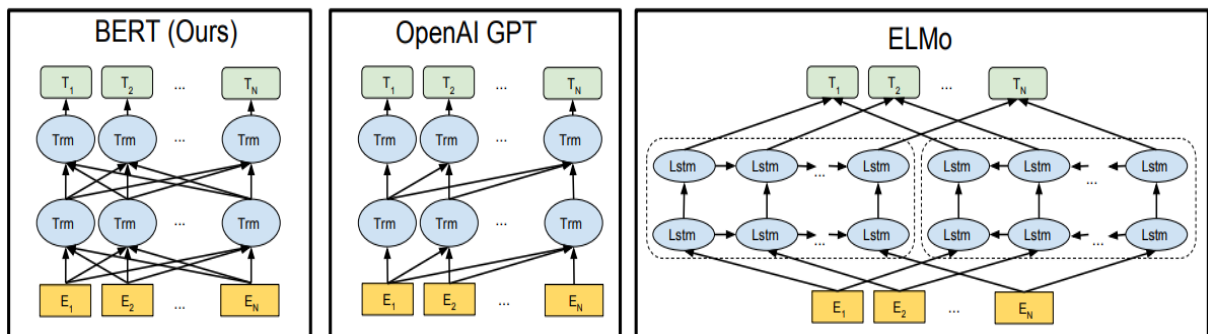


Figure 2.16: Differences between BERT, OpenAI GPT and ELMo architecture

BERT's model architecture is nearly identical to that of Transformers, specifically, a stacked bidirectional Transformer encoder without the decoder of Transformer. The differences between BERT, ELMo, and OpenAI GPT could be seen in figure 8. BERT leverage the power of bidirectional Transformer while OpenAI GPT uses a left-to-right Transformer and ELMo concatenate trained left-to-right and right-to-left LSTM.

Let L , H , and A be denoted for the number of layers (i.e. Transformer blocks), the hidden size, and the number of self-attention heads, respectively. There is two version of BERT with different sizes:

- $BERT_{BASE}$ ($L=12$, $H=768$, $A=12$, Total Parameters=110M)
- $BERT_{LARGE}$ ($L=24$, $H=1024$, $A=16$, Total Parameters=340M)

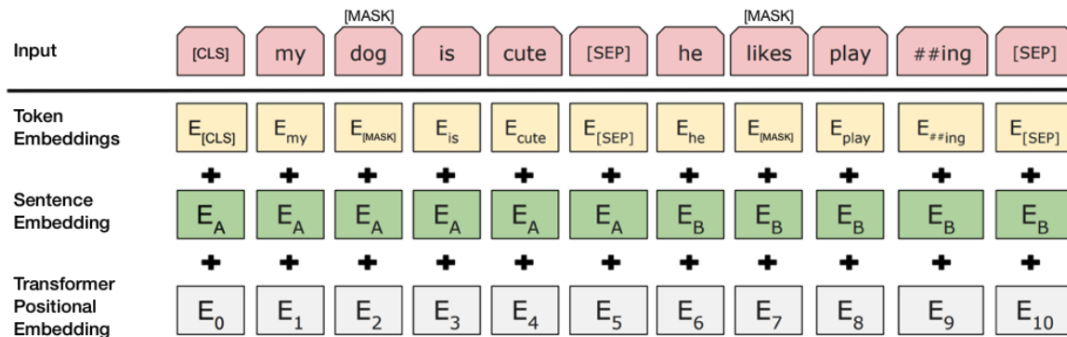


Figure 2.17: BERT's input representation

In the vanilla BERT configuration, the input representation could be one single sentence or a pair of sentences, which means the input sequence not necessarily to be an actual linguistic sentence but rather an arbitrary span of continuous text. This enables various downstream tasks (e.g., Question and Answer) to perform well with little modifications.

The initial approach of embedding is a 30000 token vocabulary WordPiece (Wu et al., 2016) [28]. The author also appends several special tokens for better sequence representation shown in figures 2.17 and 2.15.

- $[CLS]$: A special classification token stay at the beginning of the sequence. The final hidden state that maps to this token is used as a combined sequence representation for classification tasks. Therefore, this token is ignored in non-classification tasks.
- $[SEP]$: A special separate token to differentiate sentences in the input sequence. A learned embedding is also added for specifying whether each token belongs to which sentence.

Each token representation is the sum of its corresponding token, segment and position embeddings.

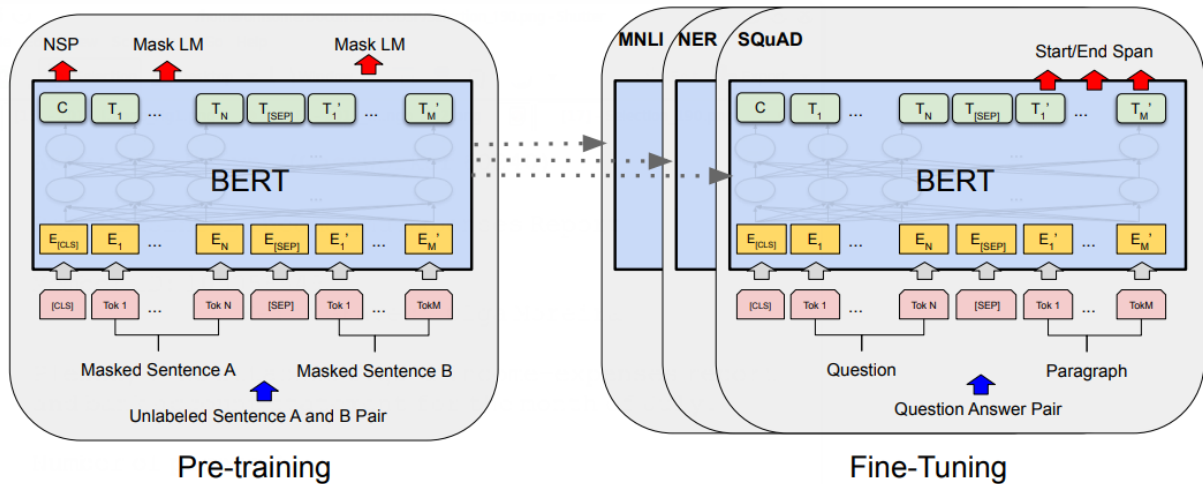


Figure 2.18: Overall pre-training and fine-tuning procedures for BERT.

2.5.3 Pre-Training BERT

Masked Language Modeling (MLM)

Language Modeling is the task of predicting the next word given a sequence of words. In masked language modeling, instead of predicting every next token, only the masked tokens are predicted. BERT is pre-trained with this task by randomly masking 15% of all tokens in a sequence and only predicting those masked words. To avoid mismatching in fine-tuning phase due to *[MASK]* token do not appear during fine-tuning, the author provides an approach that: if a token is chosen to be replaced, it will have an 80% chance to be replaced with *[MASK]* token, 10% with a random token or left unchanged for the last 10%.

Next Sentence Prediction (NSP)

Next sentence prediction task is a binary classification task in which, given a pair of sentences, it is predicted if the second sentence is the actual next sentence of the first sentence. This task is helpful for many downstream tasks such as question answering or natural language inference where they require the understanding of the relationship between two sentences. To pre-trained BERT, say choosing sentences A and B, then B would have 50% probability to actually be the next sentence that follows A (labeled as *IsNext*) while the other 50% of the time is a random sentence from the corpus (labeled *NotText*)

2.5.4 Fine-tuning BERT

The self-attention mechanism of the Transformer enables BERT to be applied in numerous downstream tasks with little modification. The input of the pre-trained BERT is similar to sentence pairs in paraphrasing, question-passage pairs in question answering, and a degenerated text- \emptyset in text classification or sequence tagging. While at the output, sequence tagging or question answering tasks would be fed in the output layer their output token representations, and *[CLS]* token is fed into the output layer for classification tasks.

2.5.5 BERT Variants

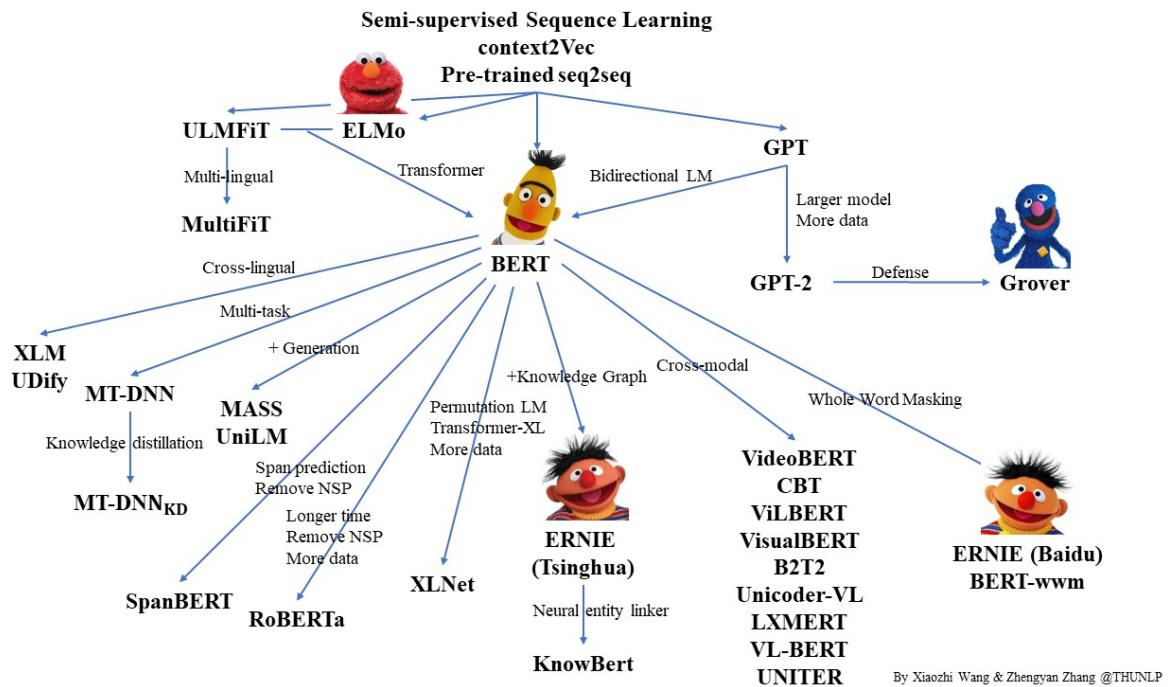


Figure 2.19: The many spawned variants of BERT

Since its release, BERT brought significant changes and has achieved state-of-the-art in multiple NLP problems, and also in other fields such as computer vision. Despite its great achievement, there are still many limitations due to its architecture and pre-trained data. Therefore, many alternative models that take BERT as its core have been proposed for many domain-specific tasks such as BioBERT (biomedical text), SciBERT (scientific publications), etc. Below are several variants that have been proved to outperform the initial BERT in many tasks:

- RoBERTa (Liu et al., 2019) [29]: The author of this model showed that the original BERT was undertrained. RoBERTa is trained for a longer period, on a bigger dataset and bigger batches with longer sequences. Especially, it also removes the task of NSP out of its pre-trained procedure and uses a dynamic masking scheme.
- XLNet (Yang et al., 2019) [30]: This model uses permutation to capture bidirectional context rather than do token masking, also combine the best of denoising autoencoding of BERT and autoregressive language modeling of Transformer-XL (Dai et al., 2019) [31].
- ALBERT (Lan et al., 2019) [32]: Known as 'A lite version of BERT', the model was proposed to enhance the training and results of BERT architecture by using parameter sharing and factorizing techniques. It reduces the number of parameters needed to train and use the task of Sentence Order Prediction rather than NSP.

Chapter 3

Related Works

3.1 Transformer-based Model for Vietnamese Spelling Correction

Research team from Vietnam National University, Hanoi proposed a deep-learning approach on the Spelling Correction task in November 2021. Their results include a correction pipeline based on transformer architecture, and a public testing set that for future papers comparison.

3.1.1 Model

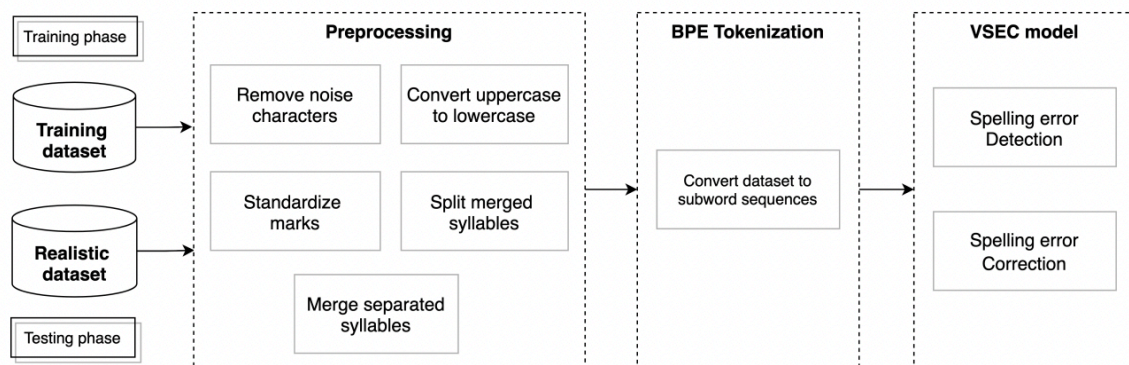


Figure 3.1: VSEC model

Preprocessing

The pre-processing phase consists of five steps:

- Step 1: Remove noise characters
- Step 2: Convert uppercase to lowercase
- Step 3: Standardize Vietnamese marks

- Step 4: Split merged syllables: Either due to user mistake or hardware error, sometimes whitespace character is not registered, which leads to merge syllables. For example: "xin lỗi" -> "xinlỗi". They utilize Peter Norvig word segmentation algorithm in conjunction to transforming into telex form to create the word.
- Step 5: Merge separated syllables: Conversely, sometimes whitespace character are misclicked, such as "chuyên" -> "chuyên". They employ the Trie structure [21], which has demonstrated its ability to browse prefixes

Tokenization

The paper opt for subword tokenization instead of character or word. The reasoning for this decision is that spelling errors comes in multiple size and shape, therefor having a syllable level will inevitably increase vocabulary size. Character level, on the other hand, increase sequence length and leads to slower training and convergence. Sub-word tokenization therefor is the perfect blend that jointly brings both contextual information of word-level and a fixed vocabulary in character level.

Model

The idea is treating Vietnamese spelling error correction as a machine translation problem. In this paper they utilize the seq-2-seq architecture provided in the paper by Vaswani et al. [20] Given an errorful sentence, we translate it into an error-free sentence. The Transformer encodes a misspelling sentence to a context hidden state using a stack of 3 encoder blocks with multi-head self-attention layer and feed-forward network. The decoder uses the encoder's hidden states and the sequence of previous target tokens to generate the target hidden state by applying a stack of 3 decoder blocks.

3.1.2 Testing set

In addition, the research team publish a public realistic testing set for comparison. The testing set were sampled from real-world sources, and hand-labeled by three different people. Contents of include 618 documents at Tailieu , an educational material website. All in all, the dataset includes 9,341 sentences, which contain 11,202 spelling errors in 4,582 different types.

3.1.3 Result and Evaluation

In general, the paper achieve positive result while using a deep-learning approach towards the spelling correction task.

Through their research, they have concluded that their approach is more effective than using a statistical based method such as N-gram. N-gram calculates the probability of a word occurring given the context, which means that it will have a hard time choosing the right word if the nearby word contains more errors. Both non-word and real-word errors can be identified and correct, however, some are limited due to domain-specific terminology.

3.2 Hierarchical Transformer Encoders for Vietnamese Spelling Correction

In May 2021, researcher from Zalo Group published a paper on Vietnamese Spelling Correction. Their paper includes a hierarchical transformer architecture that combine both word and character level information, which we will include later, and a public testing set for comparison.

3.2.1 Model

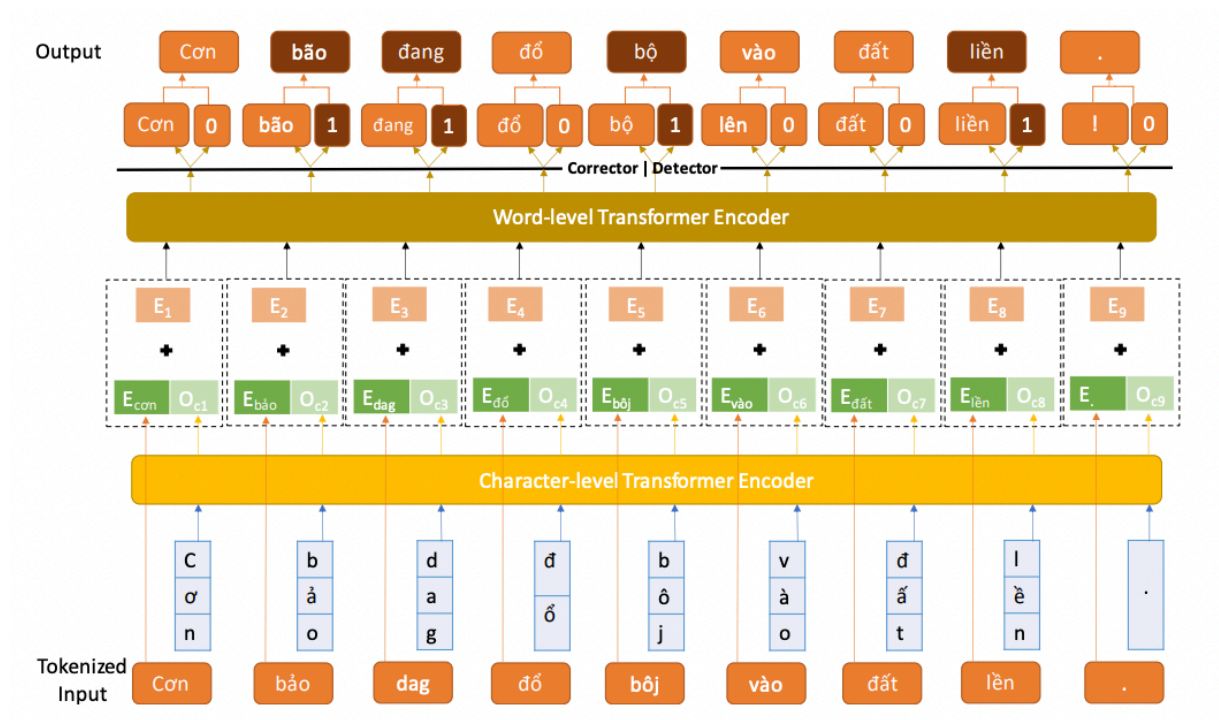


Figure 3.2: Hierarchical Transformer Model

Tokenization

Instead of opting for a hybrid subword tokenization method, the paper argues that sub-words only stand for phonetic features rather than the meaning, and that Vietnamese has a fairly small amount of tokens. Therefore, the word around is to combine both character information and word information, which contains both knowledge on the word itself and the context around it.

Model

Both character and word levels adopt Transformer encoder architectures to encode information. The character level encoder is composed of 4 self-attention layers of hidden size 256 while the word-level one is bigger, with 12 self-attention layers of hidden size 786, which is the same size as conventional BASE-size BERT. The outputs are then forwarded to two classifiers, one for detection and the other for correction.

3.2.2 Testing set

For performance measurement of the Vietnamese spelling correction problem, they collected Wikipedia drafts from varied domains. We select ones that have a substantial number of spelling errors and then manually detect the errors and suggest reasonable substitute words. In total, the test set provide more than 100 articles with about 1,500 spelling errors.

3.2.3 Result and Evaluation

The paper proposes a novel architecture based on transformer that combine both word-level information and character level information to handle the correction task, on which achieved incredible result, trained on 3Gb of clean text. In general, the paper achieve positive result while using a deep-learning approach towards the spelling correction task, while also include a public test set for the Natural Language Processing community, which we discussed in our own discussion.

3.3 Soft-masked BERT

The idea of integrating BERT into machine translation and grammatical error correction tasks is the uprising approach[6, 33, 34]. Different approaches have been implemented: Using BERT as an encoder for encoder-decoder model[33], using BERT as input embedding for a sequence to sequence model [34, 6] or using BERT as a language model.

Soft-masked BERT proposed by Zhang et al.[5] achieve impressive results for Chinese spelling correction. We would like to implement this model in Vietnamese spelling correction tasks. The model specifically tackles different sub-tasks: detecting misspelled words and correcting them based on context. The model is decomposed into two different networks specifically for those sub-tasks: a detection network and a correction network.

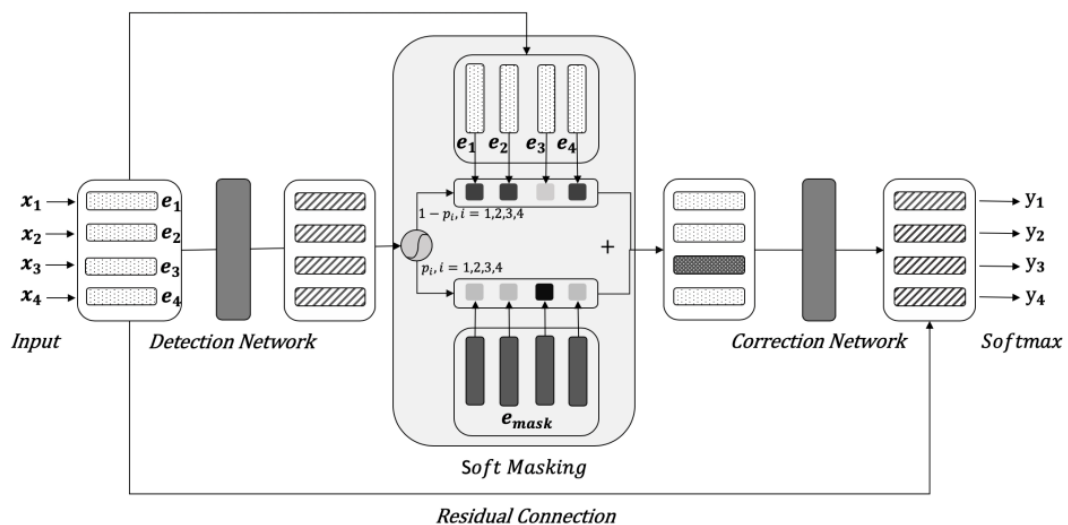


Figure 3.3: Soft-masked BERT.

The idea behind Soft-masked BERT is to re-use the already optimized pre-trained BERT model which was trained with Masked Language Modelling task. Suppose an input X is a sequence of tokens (x_1, x_2, \dots, x_n) and an output sequence $Y = (y_1, y_2, \dots, y_n)$, both sequence have the same length. If a token in X is incorrect, the model will replace that token with the correct word to obtain the correct Y . The usage of Detector and Corrector are as followed:

3.3.1 Detector

The detection network is a sequential binary labeling model, implemented as a bi-directional GRU network, where it is used to predict whether there is any token in the input sequence that is incorrect. The network input is the embedding $E = (e_1, e_2, \dots, e_n)$ from BERT embedding layer, and the output is a sequence of label $G = (g_1, g_2, \dots, g_n)$, with e_i, g_i is the embedding and label of token i^{th} , respectively. The label value of 1 indicates that token i^{th} is incorrect, and 0 if it is correct. For each token, there is a probability p_i indicates how likely that token is to be incorrect. p_i is calculated by the following equation:

$$p_i = P_d(g_i = 1|X) = \sigma(W_d h_i^d + b_d)$$

where P_d denotes the conditional probability of the detection network, σ is sigmoid function, W_d and h_d are parameters, h_i^d is the hidden state of Bi-GRU.

A soft-masked embedding e'_{mask} is created by calculating the weighted sum of input embeddings e_i and mask embeddings e_{mask} with error probabilities p_i as weights.

$$e'_{mask} = p_i \cdot e_{mask} + (1 - p_i) \cdot e_i$$

The soft-masked embedding e'_{mask} would be closer to the mask embedding e_{mask} if the probability p_i is high or else it is close to the input embedding e_i .

3.3.2 Corrector

The corrector is a BERT-based sequential model that assign probability to each of the mask embeddings $E' = (e'_1, e'_2, \dots, e'_n)$ tokens and output a sequence of tokens $Y = (y_1, y_2, \dots, y_n)$. The author take the output of the final hidden layer h_i^c and compute the probability:

$$P_c(y_i = j|X) = \text{softmax}(W(h_i^c + e_i) + b)[j]$$

where $P_c(y_i = j|X)$ is the conditional probability that token x_i is corrected as token j from the candidate list, W and b are the parameters, h_i^c is the hidden state from the last hidden layer and e_i is the input embedding of token x_i

The training of the detector and corrector is defined by the two objective functions, which are then be combined into one overall objective for learning process:

$$\mathcal{L}_d = -\sum_{i=1}^n \log P_d(g_i|X)$$

$$\mathcal{L}_c = -\sum_{i=1}^n \log P_c(y_i|X)$$

$$\mathcal{L} = \lambda \cdot \mathcal{L}_c + (1 - \lambda) \cdot \mathcal{L}_d$$

where \mathcal{L}_d and \mathcal{L}_c is the objective for training the detector and corrector, respectively, \mathcal{L} is the overall objective and $\lambda \in [0, 1]$ is coefficient of user choice.

Chapter 4

Model

4.1 Model Architecture

Inherits from the idea of VIWIKI [1] paper, instead of training from scratch the whole transformer encoder, which takes more time for the model to converge with large dataset, we use BERT as our encoder for extracting features vector from source sentence.

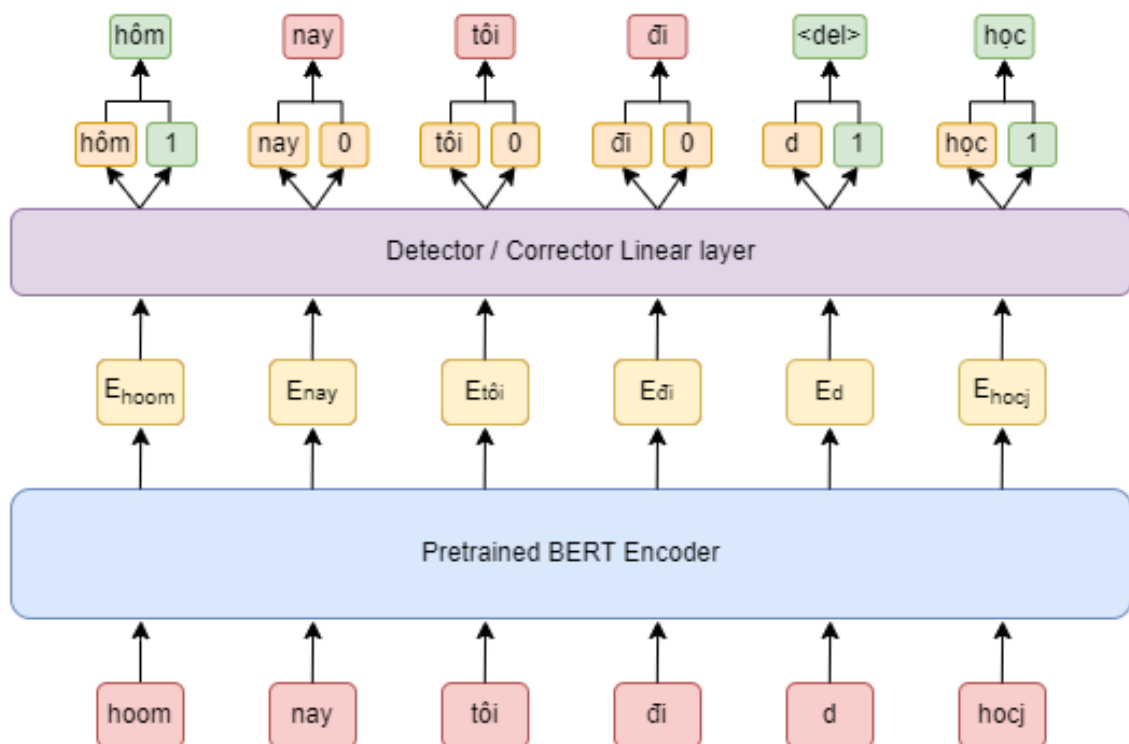


Figure 4.1: Our proposed BERT fine-tuned model

4.1.1 BERT embedding

We chose BERT as our encoder for extracting relevant information from source sentence. To be more specific, PhoBERT base setting [35] was used due to its nature of already trained on large dataset of Vietnamese text, which already captures probably all the semantics and context of Vietnamese sentences. Hence, our job is to simply add fine-tuned layers for downstream tasks.

Given a source sentence $X = \{x_1, x_2, \dots, x_n\}$, the input will go through a pretrained BERT sub-word tokenizer to achieve source input ids $E^{src} = \{e_1, e_2, \dots, e_n\}$ and attention mask $\{a_1, a_2, \dots, a_n\}$. These inputs will then be feed into BERT encoder.

$$h_{1..n}^{src} = BertEncoder(E_{x_{1..n}}^{src})$$

The encoder learn to extract relevant information of the sentence and token and output hidden representation of each token. It is then forwarded to two classifiers for detection and the other for correction.

4.1.2 Detector

In the detection phase, our purpose is to detect which token have spelling error. Hence, we cast the problem into a binary classification task, a label **0** would mean that token should be kept, else a label **1** would indicate that token is erroneous.

After getting the feature vector for BERT encoder, we take the last four hidden states of the output as our input vector for the detection phase by averaging the values in those four hidden states. A linear layer is then used for detection.

4.1.3 Corrector

After getting the output of detection phase, we forward that detection output vector into linear layer with softmax activation function. The output vector of that layer will then be concatenated with the mean feacture vector of BERT encoder that was mentioned before. The intuition behind this is to incorporate the information from the detection phase, which we believe will enhance the result of correction. Another linear layer is added for classifying the word for each tokens.

Beside classifying the tokens with the predefined BERT vocabulary, we add a new special token called $\langle del \rangle$ to indicate the deletion of token. Hence, if the corrector classify a token as $\langle del \rangle$, we would simply delete that token from the output sequence.

4.1.4 Objective function

During the training phase, cross entropy loss is used for both the detector and corrector, with an addition is that, the corrector loss will be computed using only the error token that are corrected, while the already non-error token is ignore. The final objective function would be the weighted sum of both losses:

$$\begin{aligned}\mathcal{L}_d &= -\sum_{i=0}^{N-1} \sum_{j=0}^1 y_{ij} \log p_{ij} \\ \mathcal{L}_c &= -\sum_{i=0}^{N-1} \sum_{j=0}^{V-1} y_{ij} \log p_{ij}\end{aligned}$$

$$\mathcal{L} = \mathcal{L}_d + \mathcal{L}_c$$

where N is the total number of samples, V is the vocabulary size, y_{ij} and p_{ij} stand for ground-truth labels and prediction probability.

4.2 Tokenization Repair

When fixing text errors, one particular problem arises related to the unintended white-spaces that can mess up the tokenization process. Characters that belong to another word may confuses the model as the character information is crucial to identify the correct word. Traditionally, this is a word segmentation problem, which can be solve with algorithm such as Best first search or Peter Norvig algorithm [36], segment words in a fixed vocabulary. When considering a spelling correction problem, we cannot create a vocabulary for misspelling, since it is impractical. Bast et al. [37] called this the tokenization repair problem: Given a natural language text with any combination of missing or spurious spaces, reverse it to the correct token. Fixing this type of problem is extremely hard, so the paper utilize a character language model, optionally coupled with a sequence labelling model to improve accuracy. The model calculates probability that whitespace occur after a sequence of character which we talk about below.

4.2.1 Unidirectional Language Models

The first component is a character language model, which estimate the probability of a string to occur in some language based on individual characters. The component consist of a Long short-term memory (LSTM) network that learn the context of the character stream. This is then passed through a dense layer, and a softmax output layer for character classification. The model predict $\vec{p}(s|c)$ where p is the probability that a character s comes after a context c .

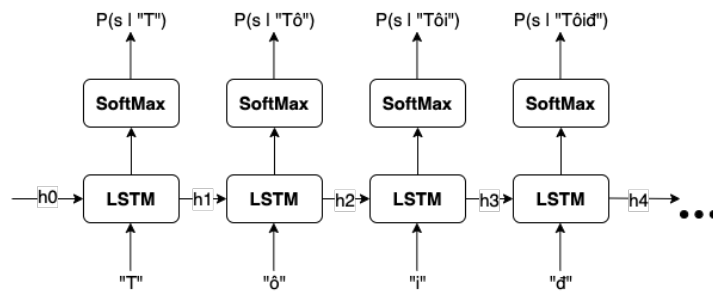


Figure 4.2: Unidirectional Language Model with LSTM

Model is trained using the categorical cross entropy loss, and sum up to 6,287,563 trainable parameters. To calculate the probability whether the space exist after a word $p(s|c)$, we define the formula:

$$p(s|c) = p(|c) * p(s|c)$$

A backward language model can easily be implemented as well, as the formula remains unchanged.

Beam search is a heuristic search algorithm that explores a graph by expanding the most promising node in a set limited by parameter called **beam** or beam-width. At each step of the search, algorithm keep track of the score and the partial solution string in the search state.

At each timestep i , the search state is determined by whether a space is added before adding the character T_i or not. Without adding space, search state is:

$$S_i = S_{i-1} - \log(\vec{p}(T_i|R_{i-1})) + P_{del}$$

$$R_i = R_{i-1}T_i$$

With adding space, search state is:

$$S_i = S_{i-1} - \log(\vec{p}(T_i|R_{i-1})) + P_{ins}$$

$$R_i = R_{i-1}T_i$$

Notice that we add 2 penalty score for inserting and deleting a space. However, this penalty is only added if the change does not exist in the original sentence, otherwise, it is zero.

4.2.2 Bidirectional Sequence Labeling Model

The second component is a sequence labeling model. Given a sequence of non-space string of character from the front and back, the model learn to predict the probability \overleftrightarrow{p} of a space occur in between. Here we again use the same RNN model as the language model, only now utilizing a bidirectional version. The model is trained with binary cross entropy loss.

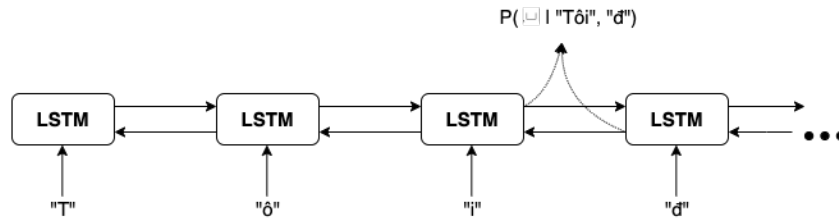


Figure 4.3: Bidirectional Sequence Labelling with BiLSTM

To combine both model, we can slightly modify beam search. At each timestep i , the search state is determined by whether a space is added before adding the character T_i or not. Without adding space, search state is:

$$S_i = S_{i-1} - \log(\vec{p}(T_i|R_{i-1}) * (1 - \overleftrightarrow{p})) + P_{del}$$

$$R_i = R_{i-1}T_i$$

With adding space, search state is:

$$S_i = S_{i-1} - \log(\vec{p}(T_i|R_{i-1}) * \overleftrightarrow{p}) + P_{ins}$$

$$R_i = R_{i-1}T_i$$

4.3 Pipeline

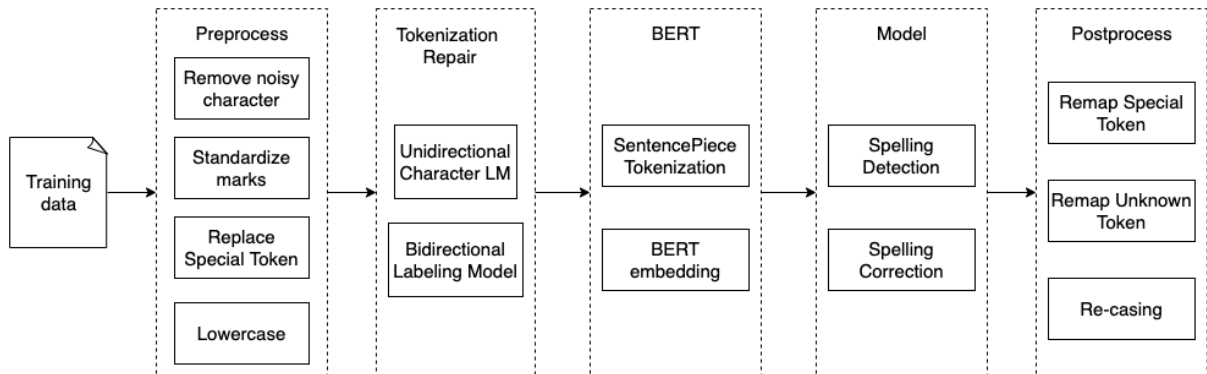


Figure 4.4: Spelling correction pipeline

Given a sentence with spelling error, the pipeline is as follow: **Preprocess, Tokenization, Model inference, Postprocess**. For preprocess:

- **Remove noisy character**: clean unnecessary character such as emoji, leading spaces, etc.
- **Standardize mark**: the data is standardized using a mapping set between the telex syllable and the correct mark syllable. For example: “cuả” -> ”củ”
- **Lowercase**: Keeping original casing with significantly increase vocabulary. Therefore, we lowercase token for smaller and faster model.
- **Tokenization repair**: Sentence pass through the aforementioned tokenization repair model to fix whitespace error, and fix merged tokens.
- **Replace Special Token**: some token like numbers and dates can increase vocabulary, all the while not adding context to the sentence, hence we can replace it with one special token like *date* or *number*.

In the tokenization phase, we use the default subword tokenizer used by PhoBert, which has the vocabulary size of 64000. After that, the tokens are fed to BERT model to retrieve embedding. average the sub-word representations to obtain the word representations. This is then fed to 2 linear layer, detection and correction. Detection layer produce the probability of a token begin misspelled. Then, correction layer will perform token classification to predict the correct word. Finally, we perform post-process on the model output:

- **Remap Special Token**: replace special token like dates and number with the original token we have saved in preprocessing step.
- **Remap Unknown Token**: replace unknown token by aligning output with source sentence, then map all unknown token with original token.
- **Recasing**: perform recasing for special cases like Capitalize, CAPLOCKS.

Let us take an example of a sequence-to-sequence model using sub-word tokenization:

Input Text	"Ngafy 14/3, học sinh toànquosc được nghỉ học."
Pre-processing	"ngafy date học sinh toàn quosc được nghỉ học"
Tokenization	["nga@@", "f@@", "y", "date", "học", "sinh", "toàn", "quosc", "được", "nghỉ", "học"]
Model Inference	["ngày", "date", "học", "sinh", "toàn", "quốc", "được", "nghỉ", "học"]
Post-processing	["Ngày", "14/3", "học", "sinh", "toàn", "quốc", "được", "nghỉ", "học", "."]
Output Text	"Ngày 14/3, học sinh toàn quốc được nghỉ học."

Table 4.1: Example of pipeline

Chapter 5

Application Implementation

In this section, we would briefly introduce several tools and technology aspect of our implementation while building these models, namely, programming languages, frameworks and public tools.

5.1 Programming Languages

5.1.1 Python

Python is a high-level, interpreted, general-purpose programming language created by Guido Van Rossum and firstly introduced in 1991. Its design philosophy emphasizes code readability with many advantages, specifically:

- Simplified syntax with emphasis on natural language, which make it easier for beginners to learn and also enhance readability.
- Its enormous and ever-growing ecosystem of libraries and packages. Python could be used for developing any kind of applications in any industry or field.
- It is famous for solving data analysis problems and building machine learning and deep learning models with the provision of many powerful mathematics and visualization framework such as TensorFlow, PyTorch, Pandas, NumPy, Scikit-learn, etc.

In this project, we utilize **PyTorch** as our main library for building models due to its support for coding machine learning modules in object oriented ways which also enhance readability. **Transformers** libraries from HuggingFace [38] is also used for obtaining BERT's pretrained model and tokenizer

For back-end development, we use FastAPI for building APIs for our application. **FastAPI** is a modern, fast with high-performance web framework for building APIs with Python. Is is claimed to be one of the fastest Python frameworks available with the support of handling asynchronous requests and performance on par with Go and NodeJS while being lighter weight comparing to exist Python API web frameworks such as Django. Built-in support such as Pydantic and SwaggerUI are also some of its main advantages for making APIs easier to maintain and deploy.

5.1.2 JavaScript

JavaScript, developed by Brendan Eich at Netscape in 1995, first appeared as Mocha, then to LiveScript and lastly JavaScript, is a programming language that have its syntax similar to Java. It is mostly used for developing front-end of applications, for both websites and mobiles, with its main functionality is to navigating front-end workflows.

We use JavaScript with ReactJS as our tool for developing the front-end for our demonstration website, with the simple use case of inputting a text field for finding errors and correcting it.

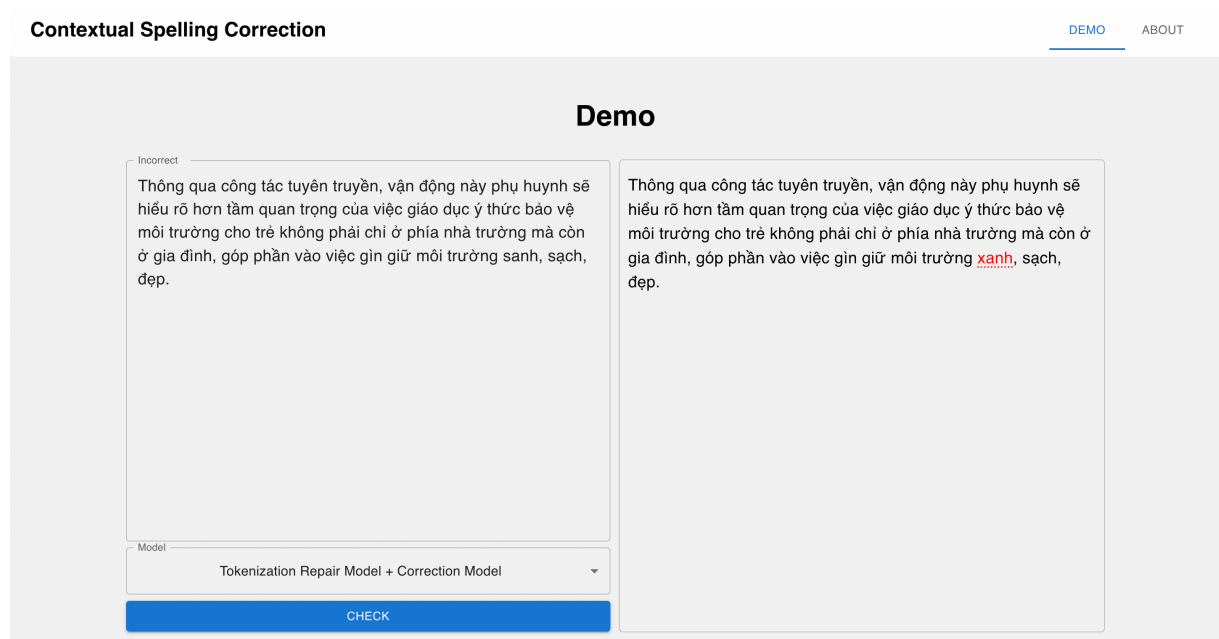


Figure 5.1: Web-app interface to demonstrate result.

Demo

Incorrect

Thông qua công tác tuyên truyền, vận động này phụ huynh sẽ hiểu rõ hơn tầm quan trọng của việc giáo dục ý thức bảo vệ môi trường cho trẻ không phải chỉ ở phía nhà trường mà còn ở gia đình, góp phần vào việc gìn giữ môi trường sanh, sạch, đẹp.

Thông qua công tác tuyên truyền, vận động này phụ huynh sẽ hiểu rõ hơn tầm quan trọng của việc giáo dục ý thức bảo vệ môi trường cho trẻ không phải chỉ ở phía nhà trường mà còn ở gia đình, góp phần vào việc gìn giữ môi trường xanh, sạch, đẹp.

Correction Model

Tokenization Repair Model

Tokenization Repair Model + Correction Model

Figure 5.2: Choose a model to demonstrate performance

5.2 Tools

5.2.1 Google Colaboratory

Google Colaboratory is a free-to-use cloud service for analysing data and training machine learning models. Google colaboratory inherits the interfaces and functionalities from Jupyter Notebook, which can only be run on local machine, with the extension of storing data with Google Drive and providing free GPU (mostly Tesla K80) with a fully-equipped machine learning virtual environment.

Chapter 6

Result and Evaluation

6.1 Dataset

For the spelling correction task, there are no common data for us to train and compare with other model. Fortunately, Vietnamese spelling error follows some common rules that we can implement rule-based random augementer. Through some research, we determined the rules to augment an error-free text into an errorful text.

6.1.1 Preprocessing Data

Our pre-processing pipeline goes through the steps below:

- Step 1: Remove noise characters: We clean the following illegal characters before augmentation:
 - Compound unicode to standard unicode
 - Remove links
 - Remove HTML tags from crawling
 - Remove emoji
 - Remove trailling space and newlines
- Step 2: Convert all characters to lowercase: Keeping it uppercase will drastically increase vocabulary size, which decrease the model performance.
- Step 3: Separate punctuation mark: Like uppercase, having an additional punctuation mark to the word does not add any information.

6.1.2 Errors and Augmenters

Errors Definitions

Firstly, we have to understand different types of misspellings. In the spelling correction task, spelling can be divided into two kinds: non-word errors and real-world errors[39].

In layman's term, non-word errors are misspelled word create an out-of-vocab word; for example, "trừ" may be misspelled to "truwf", due to a telex mechanism in Vietnamese keyboard. Accidental insertion and deletion of characters also fall into this category, which we will discuss

later.

Real-world, on the other hand, generates other word inside the vocabulary, which makes it harder to detect and model must rely on context to correct it. Vietnamese spoken language is diverse and varies from region to region, an instance of this would be "vui vẻ" and "dui dễ", where both words phonetically sound the same but are spelled differently. Shortcuts and teencodes are not considered real-world errors.

If we wish to generate data that will reflect real-world data. To have a grand view of different types of errors, we firstly crawl a set of data from multiple sources, including 'vi.wikipedia.org', lyrics and movie transcript, and from Vietnamese literature.

Different error type in Vietnamese

- **Typos:** Vietnamese typists mainly use telex or VNI typing mechanism for pure characters (á, ã, â). These characters are compounded by multiple character, for example, á is a combination of a and s in telex. Random errors by typing text such as character insertion, deletion, duplicate, swap are also common when typing fast, or close-proximity character on the keyboard. Inappropriate white-space insert or deletion are not needed as the tokenization repair model handles white-space restoration.
- **Spelling Errors:** In different regions, people may use unofficial or regional dialect. This may include d/r/gi, tr/ch, s/x, n/l mixture at beginning of sentence, and so on.
- **Diacritic-related error:** diacritics are an important part of the Vietnamese alphabet. In the written language, vowels (a, e, i, o, u) can be comprised of more than just the base but also diacritics to form a full word. When a sentence is full of missing diacritics, it is harder to deduce the meaning of a word. Numerous scholars have studied the field of Vietnamese Diacritics Restoration problem [40, 41, 42] with positive results. Non-word errors in this category comprise of fully missing diacritics or partially missing diacritics, which Pham et al. [41] has classified as type I error. These errors can be predictable and can be easily regenerated.

Index	Error	Example
1	Telex Error	"uống" => "uoosng", "uoongs", .. "trăng" => "trawng", "trangw", ..
2	VNI Error	"uống" => "uo16ng", "uong61", .. "trăng" => "tra7ng"
3	No dialect	"ăn uống" => "an uong"
4	Missing dialect	"một ngày" => "môt ngày"
5	Spelling errors	"không" => "khôn", "khổng", "hông", .. "thôi" => "thô", "hôi", "thi", ..
6	Start with d/r/gi/v	"dang rộng" => "vang rộng"
7	Ends with c/t	"cá cược" => "cá cươt"

Table 6.1: Example of some misspelling errors

6.1.3 Augmentation Pipeline

Our augment process guarantee a more uniform distribution of errors across all token in the vocabulary. First, we build a vocabulary from the corpus and mark the sentence that contains the word within it. For each words, we sample a specific number of marked sentence for each

type of error and replace error token based on the augementer output. We also maintain the error percentage of a sentence to be lower than 30%, so the marked sentence list has to be re-evaluate before sampling. This process is repeated until all word in the corpus has been augmented.

In our augmentation pipeline in the thesis proposal phase, we randomly augmented a sentences with a random augementer. Each sentence has 30% error, and the augementer is randomly chosen. Compared to this method, our current method guarantee that each word will equally be augmented with the same amount of error percentage, all the while maintaining the error percentage of a sentence to a certain percent. We also have finer control over which word will be augmented, specifically Vietnamese words.

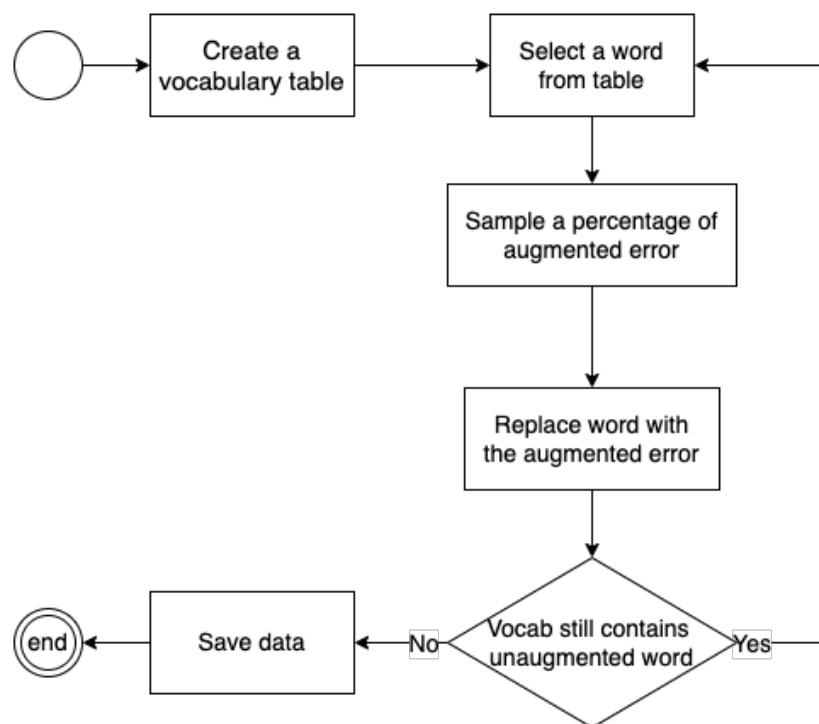


Figure 6.1: Data augmentation pipeline

6.1.4 Overall Data

Our testing data contains a little over 5000 sentences with various errors mentioned above. However, depending merely on our testing set is insufficient to evaluate the accuracy in real-life tests. Therefore, we include the testing dataset made public by [9] and [1].

In total, we achieve the following

Dataset	Pair of sentences
Training set	2.518.213
Validation set	171.910
Testing set	5.000
VSEC set	9.341
Viwiki set	14.831

Table 6.2: Data for spelling correction

6.2 Real world data analysis

To observe on the frequency of error types, we crawl several small corpus from informal news websites such as kenh14.vn, tintinonline.com.vn, etc.

In total, we achieve the following statistics:

	Counts
Number of sentences	1067
Total tokens	48816
Total error tokens	573

Table 6.3: Real world data counts

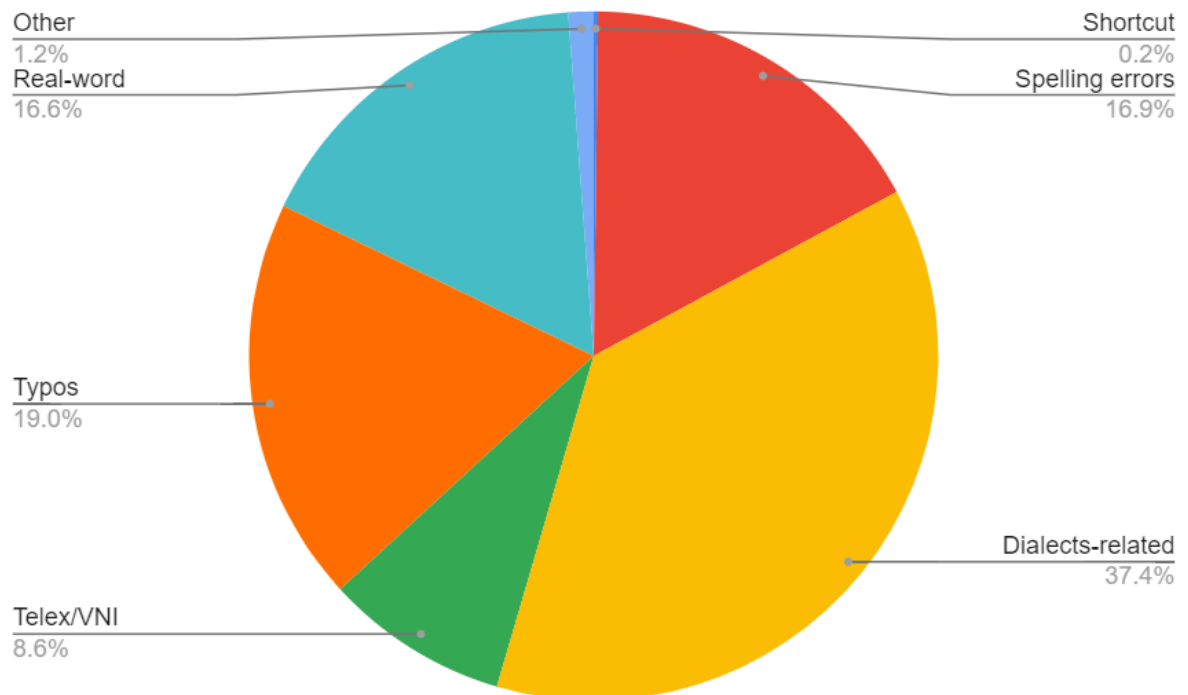


Figure 6.2: Errors type percentages

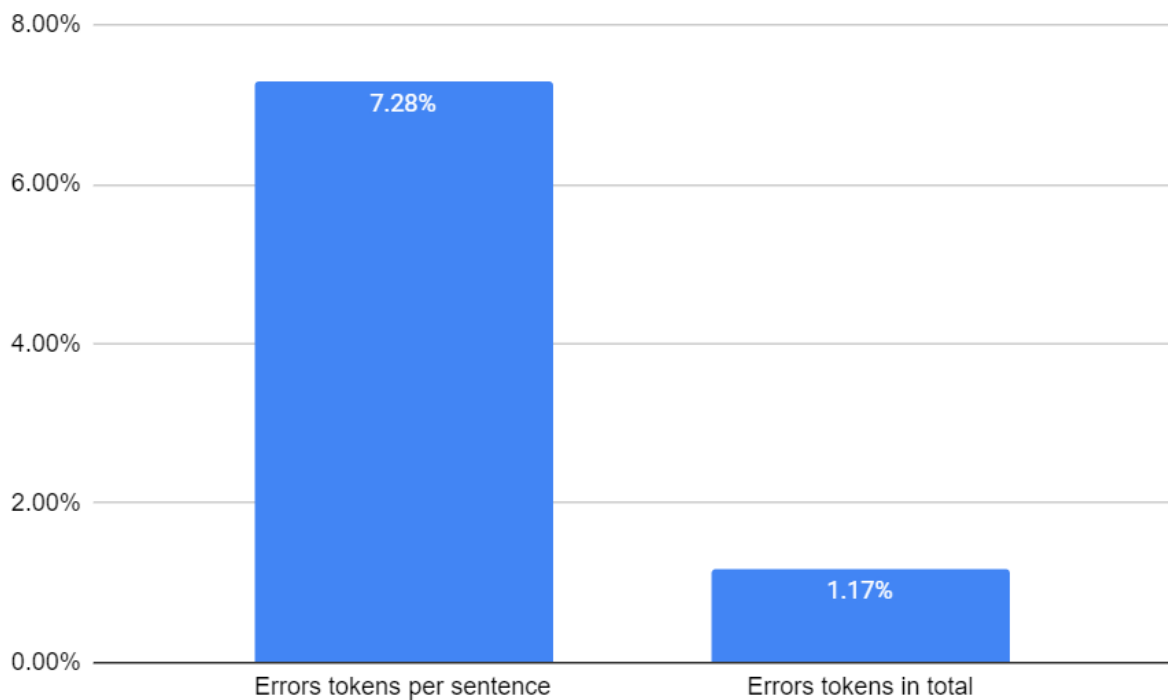


Figure 6.3: Percentages of error tokens

While the number of error tokens only takes up at just above 1% of total tokens, on sentence that has errors, the percentage of errors tokens amount for over 7%. With most of the errors are related to dialect mark, such as missing or wrong dialects. Typos is another common error where having keyboard mis-clicked often occurred. Not to mention, spelling errors also has high rate of occurrence, words with missing characters is the most common type.

6.3 Baselines and Model

6.3.1 Baselines models

We have chosen several models as the baseline for this thesis. The first one is the transformer architecture proposed by Vaswani et al.[20]. For the transformer model, we additionally train on different tokenization methods to compare the effectiveness. We also experiment with Soft-mask BERT [5], which has achieved near the state-of-the-art result for Chinese, and a modified version of its to study the efficiency of pretrained models. Finally, we use BERT Fine-tuned with two layer of correction and detection as our model for further evaluation and has achieve significant results.

6.3.1.1 Transformer

The Tranformer model [43] was trained also using the implementation from Fairseq. We decide to use the BASE model proposed in the paper. To compare the effectiveness of different type of tokenizer, we decide to train all word, subword, character level tokenization. For the subword tokenization, we limit the vocabulary size to 16000.

6.3.1.2 Soft-masked BERT

Soft-masked BERT is a model proposed in Chinese Spelling Correction. The model reach State-of-the-art result for Chinese, so we decide to include the model for Vietnamese.

6.3.1.3 Hard-masked XLMR

In the process of training the Soft-masked BERT model, by trial and error, we find out that in order to make this model perform well, an extensive amount of training time is needed due to its heavy architecture and large number of parameters (550M parameters).

With the same approach as Soft-masked BERT, we re-use the architecture of the model, called Hard-masked BERT, with the only difference is that instead of using BERT to learn and correct the error, we use a pre-trained model that was trained with Masked Language Modelling for the corrector.

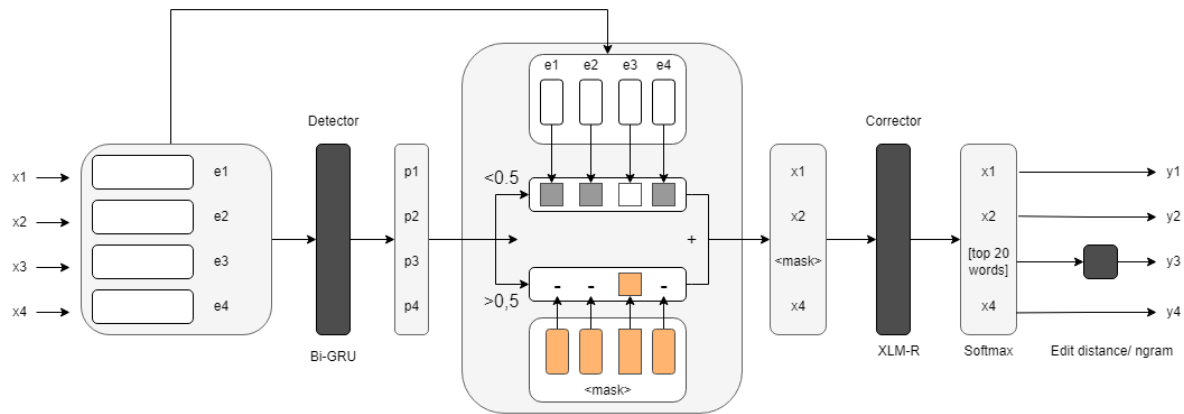


Figure 6.4: Hard-masked BERT.

Detector

From figure above, the detector of Hard-masked BERT is the same as the soft-masked version, a Bi-GRU network. While the input of the detection network in soft-masked is the embeddings from BERT, we instead use a subword tokenizer trained with the SentencePiece package for ease of use. The output is also a sequence of labels where we use a sigmoid function to calculate the probability of the error tokens and use a threshold of 0.5 to manually mask the tokens.

Corrector

We chose pre-trained XLM-Roberta (Conneau et al., 2019) [44] as our corrector. It is a cross-lingual masked language model based on RoBERTa and has been pre-trained on 100 different languages, including Vietnamese.

From the detector, a sentence with mask tokens is then processed to align with the model input settings, then pass these input into the model. A softmax function is used to find the probabilities for the candidate tokens so that we could choose 20 tokens with the highest probabilities.

For the task of choosing the appropriate token, we compute the Levenshtein similarity ratio

of the incorrect token with each of the candidate tokens and eliminate all the tokens with a ratio less than 0.4. By doing this, we could assure that the candidate token is not completely different from that of the error token. Although, in many cases, several candidate tokens will have the same edit distance. Hence, a tri-gram model is used to score whether which token is suitable in the context of the input sequence and then pick the token with the highest score. Another version which post-processing is not applied but instead use the best candidates token via probabilities from the corrector will also be tested for comparison.

This approach, however, will not always correct the spelling mistake but rather choose the word that best suits the context based on its learning phase.

6.3.2 Configuration settings

6.3.2.1 Transformers

For all three different type of tokenization of transformes, the configuration is as below:

- Encoder layer: 3
- Decoder layer: 3
- Encoder attention head: 8
- Decoder attention head: 8
- Encoder embedding dimension: 512
- Decoder embedding dimension: 512
- Encoder feed-forward dimension: 2048
- Decoder feed-forward dimension: 2048
- Drop out: 0.1
- Learning rate: $1 * 10^{-4}$
- Learning rate scheduler: inverse square root
- Batch size: 128

6.3.2.2 Soft-masked BERT and Hard-masked XLMR

For both Soft-masked BERT and Hard-masked XLMR, we use the following configuration settings:

- Learning rate: $1 * 10^{-5}$
- RNN layer: 2
- Feed forward dimension: 512
- Batch size: 32

6.3.2.3 BERT Fine-tuned

- Learning rate: $1 * 10^{-5}$
- Feed forward dimension: 512
- Batch size: 64
- Vocab size: 40000

6.3.3 Overall training settings

The training process requires a significant amount of resources like RAM and GPU. Therefore, the training and inference process is done on the Google Colab Pro environment, which provides us with resource as listed:

- Operation System: Linux
- RAM: 12Gb
- GPU: Nvidia Tesla P100

All models are trained with Adam optimizer [45]. In addition, we also utilize mixed precision training [46] to boost model training time.

6.4 Evaluation metrics

6.4.1 Detection and Correction

We discuss evaluation methods to measure how well the models perform. Spelling correction problem can be viewed as a binary classification of prediction, where the matching elements are considered as true predictions and the others, incorrect predictions.

In the past, studies for Vietnamese spelling correction often opt for word-level evaluation of precision, recall, and f-score [8, 7]. These metrics can be applied to both detection and correction tasks. In this work, we will use both the metrics that are provided by VSEC [9] and VIWIKI [1] paper.

$$DetectionPrecision = \frac{TrueDetections}{TotalErrorDetected} \quad (DP)$$

$$DetectionRecall = \frac{TrueDetections}{TotalActualErrors} \quad (DR)$$

$$DetectionF1 = \frac{2*DP*DR}{DP+DR} \quad (DF)$$

$$CorrectionPrecision = \frac{TrueCorrections}{TotalErrorDetected} \quad (CP)$$

$$CorrectionRecall = \frac{TrueCorrections}{TotalActualError} \quad (CR)$$

$$CorrectionF1 = \frac{2*CP*CR}{CP+CR} \quad (CF)$$

$$Accuracy\%_{detected} = \frac{TrueCorrections}{TrueCorrections+WrongCorrections} \quad (Acc_d)$$

$$Accuracy_{total} = \frac{TrueCorrections}{TrueCorrections+WrongCorrections+WrongDetections} \quad (Acc_t)$$

6.4.2 Tokenization Repair

Given a corrupt input text I, a ground truth text T and a predicted text P, the algorithm predicts a set of space insertions and deletions that aims to transform C into T. The paper use two metrics for the evaluation: F-score and sequence accuracy.

F-score: We define $edits(A, B)$ as the space insertions and deletions that transform A into B. To achieve alignment on both sequence, we can use levenstein dynamic alignment algorithm on character level. If we let $C = edits(C, T)$ be the ground truth edit operations and $P = edits(C, P)$ the predicted edit operations, the number of true positives is $TP = |C \cap P|$, the number of false positives is $FP = |P \setminus C|$ and the number of false negatives is $FN = |C \setminus P|$. The F-score is computed as:

$$F - score = \frac{TP}{TP + \frac{1}{2} * (FP + FN)}$$

Sequence accuracy: sequence accuracy measures what fraction of exact correction ($P = T$) over the whole testing set.

6.5 Result And Evaluation

We tested our models against three test set. Firstly, we augmented a small test set of 5000 sentence that has not been seen in the training set. Our augmentation includes all errors that are included in the training set.

Models	Detect			Correct				
	DP	DR	DF	CP	CR	CF	Acc_t	Acc_d
Transformers-Char	96.22	73.24	83.17	71.73	54.59	62.00	71.73	74.54
Transformers-Word	71.89	84.23	77.57	40.00	46.86	43.16	40.00	55.63
Transformers-Subword	98.84	85.17	91.50	85.29	73.49	78.95	85.29	86.29
Hard-Masked XLMR	97.82	93.26	96.48	46.33	55.6	50.54	46.33	53.66
Soft-Masked BERT	96.36	93.92	95.12	72.35	68.73	70.49	72.35	75.12
Our Model	97.20	<u>96.42</u>	<u>95.62</u>	<u>91.70</u>	<u>89.79</u>	<u>90.73</u>	<u>91.70</u>	<u>93.12</u>
+ Tokenization Repair	<u>98.36</u>	97.31	97.83	92.07	91.09	91.58	92.07	93.61

Table 6.4: Results

Out of the baseline transformers, we can see that subword tokenization clearly outperforms other character and word level model at 91.50% in terms of detection. We predicted the word-level model performs the worst since multiple non-word errors can be generated from one correct word, hence our dataset can not cover all cases of different errors. Character tokenization does perform better than word tokenization, although we note that inference speed of both word and character model performs more poorly, especially the latter. In terms of correction, subword-tokenization also achieve notable result, where its accuracy surpass other transformers baselines.

Soft-masked BERT model (which achieved SOTA in Chinese Spelling Correction) also performs poorly as it also uses word-level tokenization, which leads to a very large vocabulary and slower convergence rate. Its modified version, Hard-masked XLMR, does not achieve any higher result than the base model since it only correct the mask tokens based on the surrounding context rather than the error token information. However, both the model has improve the detection phase by incorporating a detection layer.

Our main model detection achieve f-score around 95.62%, which outperform other baselines. The results further improve as we add the tokenization repair module which increase to 97.83%. In terms of correction, we achieve 93.61% in accuracy within the detected token and 92.07% accuracy in general. However, in real life testing, our model still meets some limitation such as teencode, shortcut terms, and foreign word which we will touch in the discussion.

	Sentence
Incorrect	Thông qua công tác tuyên truyền, vận động này phụ huynh sẽ hiểu rõ hơn tầm quan trọng của việc giáo dục ý thức bảo vệ môi trường cho trẻ không phải chỉ ở phía nhà trường mà còn ở gia đình, góp phần vào việc gìn giữ môi trường sanh , sạch, đẹp.
Correct	Thông qua công tác tuyên truyền, vận động này phụ huynh sẽ hiểu rõ hơn tầm quan trọng của việc giáo dục ý thức bảo vệ môi trường cho trẻ không phải chỉ ở phía nhà trường mà còn ở gia đình, góp phần vào việc gìn giữ môi trường xanh , sạch, đẹp.
Predict	Thông qua công tác tuyên truyền, vận động này phụ huynh sẽ hiểu rõ hơn tầm quan trọng của việc giáo dục ý thức bảo vệ môi trường cho trẻ không phải chỉ ở phía nhà trường mà còn ở gia đình, góp phần vào việc gìn giữ môi trường xanh , sạch, đẹp.
Incorrect	tôiddens đây v ới tu cách là người phán xử.
Predict	tôi đến đây với tu cách là người phán xử.

Table 6.5: Example of correction

Next, we evaluate our model on VSec public test set by Vietnam National Univerisity team.

Models	Detect			Correct				
	DP	DR	DF	CP	CR	CF	Acc_t	Acc_d
Transformers-Char	37.39	56.61	45.03	26.13	39.57	31.48	26.14	69.91
Transformers-Word	21.69	66.68	32.73	14.64	45.01	22.10	14.64	67.51
Transformers-Subword	82.00	76.95	79.39	71.82	67.40	69.56	71.82	87.59
Hard-Masked XLMR	85.03	78.11	82.24	53.91	42.44	47.49	53.91	68.78
Soft-Masked BERT	81.55	79.2	80.35	59.8	56.76	58.24	59.8	78.43
Our Model	92.20	<u>85.23</u>	<u>88.58</u>	83.70	<u>77.36</u>	80.40	83.69	90.77
+ Tokenization Repair	93.59	85.49	89.36	86.23	78.76	82.33	86.23	92.12
VSEC (2M)	89.1	76.9	82.6	82.6	71.3	76.5	-	-
VSEC (5M)	93.1	81.3	86.8	87.4	76.3	<u>81.5</u>	-	-

Table 6.6: VSEC test set

Bold for highest score, underline for second highest

Each sentence in the VSEC test set is segmented and contains at least one error. As an overview, the set include all of our defined error, plus some minor cases of adding token, which our model fails to correct due to the fact that we only predict whether to keep or change the corresponding token. The table above includes VSEC model trained on two million sentence dataset and five million sentence dataset. In the detection task, our model overall performs better than the paper best result. Precision score is lower by 0.6%, while our recall score is higher by 4.76%. Overall, the f-score with our model increased by 2.36%. When compared with the VSEC model trained with two million sentence, our model achieve higher score in all three benchmarks.

For the correction task, although our model perform worse on correction by 2.37%, we outperform on recall performance and f-score performance, by 2.81% and 0.47%, respectfully. Once again our model outperform all score compared to the VSEC model on two million dataset.

Models	Detect			Correct	
	DP	DR	DF	Acc_t	Acc_d
Transformers-Char	1.33	49.08	2.6	0.65	48.93
Transformers-Word	1.08	67.17	2.14	0.46	42.81
Transformers-Subword	4.38	58.64	8.16	3.36	76.73
Hard-Masked XLMR	23.84	60.67	34.22	8.96	48.32
Soft-Masked BERT	21.72	58.76	31.72	12.38	60.54
Our Model	40.26	63.49	49.27	35.84	89.03
+ Tokenization Repair	<u>40.65</u>	<u>64.87</u>	<u>49.98</u>	36.29	<u>89.27</u>
Viwiki	66.96	70.92	68.88	64.29	96.01

Table 6.7: Viwiki test set

Bold for highest score, underline for second highest

Viwiki dataset comprised of all a hand-labeled text crawled from Vietnamese wikipedia. Unfortunately, our model performs much worse compared to the paper model. We offer some detailed explanation in the discussion. In general, test cases in this test set does not only includes Vietnamese spelling error, but also English spelling errors, which tremendously reduce our accuracy. All the model of our choices have high recall score and low precision score in the detection phase, which implies that false positive samples are significant to that of false negative. In other words, our model predict a large portion of tokens as erroneous while in fact they are not. These wrong detections also affect our accuracy of correction in total. Although, the accuracy in detected portion is 88%, which means, for every token that are actual errors, our model have an average accuracy of 88% of correcting them to true tokens.

We also measure our tokenization repair module, which achieve a high score on our augmented test set. Unfortunately, we could not find a real test set in Vietnamese to evaluate.

Model	F-score	Sequence Accuracy
UNI	93.7	97.1
BID	97.5	98.1

Table 6.8: Tokenization repair result

UNI: Unidirectional Language Model, BID: Bidirectional Label

We take a few example to demonstrate our results.

	Sentence
Incorrect	Khôn gxar rác để giữ một môi rường xanh, xạch, đẹp
Predict	Không xar rác để giữ một môi rường xanh, xạch, đẹp
Incorrect	tôiddens đây v ới tu cách là người phán xử.
Predict	tôi ddens đây với tu cách là người phán xử.

Table 6.9: Example of tokenization repair

As we can see, the tokenization repair module are able to correctly predict the space insertion and deletion in the presents of spelling mistakes. This is a crucial improvement for our fine-tuned model as it does not have the ability to seperate merged syllable, such as "tôiddens", or merge splited tokens such as "v" and "ới".

	Sentence
Original	qua công tác tuyềntryuen, phụ huynh cần phain ắm rõ được tầm quan trọng trong việc nhận thức của con.
Without TR	qua công tác tuyềntryuen, phụ huynh cần phải nắm rõ được tầm quan trọng trong việc nhận thức của con.
With TR	qua công tác tuyên truyền, phụ huynh cần phải nắm rõ được tầm quan trọng trong việc nhận thức của con.
Original	tôi dduow c đi học
Without TR	tôi đưa c đi học
With TR	tôi được đi học

Table 6.10: Example of correction with tokenization repair

TR: tokenization repair

6.6 Effect of Hyperparameter

We also investigate further the effect of larger dataset on the accuracy of the proposed method. Since we do have a limited resource, our training set max out at around 2.500.000 sentences. The result is displayed in the table below.

Training Set	Detect			Correct	
	DP	DR	DF	Acc_t	Acc_d
500k	90.37	89.98	90.17	85.70	90.60
1M	92.98	92.40	92.69	87.13	91.24
2M	94.50	93.76	94.13	89.91	92.08
2.5M	98.36	97.31	97.83	92.07	93.61

Table 6.11: Effect of training set size on accuracy

Chapter 7

Discussion and Future Work

Firstly, we will discuss on why we leave out several common errors such as teen-code, shortcut terms and also the problem of dealing with foreign words.

Teencode

While teen-code is quite commonly used in daily text and messages, it is an intentional attempt by user to stylise the word, hence does not fall into the spelling correction category. Also, we could not recreate this type of error without spending hours labeling them. Due to its uniqueness and ruleless, dealing with them is painful.

For instance, someone may put in a number “3” to replace the letter “e”, randomly drop some letters, or add random punctuation in-between letters and words. A particular teen-code sentence would be as follow: “b4.nh dzạo n3‘ thja’ nà0 r04j” which translates to “Bạn dạo này thế nào rồi”.

Moreover, the teen-code “dictionary” tends to increase drastically, with different ways of creating them are being “invented” almost daily and, not to mention, the emergence of trending words. Human may not understand every words written with teen-codes, even given the whole context of the passage. Hence, correcting them is not always possible.

Shortcut terms

As for shortcut terms, we could simply replicate this error by only taking the first letter of each words and then concatenate all those letters. However, not every words could be shorcuted. These shortcut terms tends to be relevant to one single particular domain. While several terms could be used in general purpose, these terms might be abbreviated to various meaning depending on the context.

A simple example of the term “bn” could stands for ”bạn” or “bao nhiêu”. In the context of clothing shop, one may text “Bộ này bn vậy bn?”, which could get us confused the words ordering.

In order to solve this, building a specific dictionary for each fields is a must if we want to augment this type of errors. Again, time and resource limitation restrict us from doing so.

Foreign words

Although foreign words may still appear in text and convey meaning, we decided to exclude

them from our correct scope for a few reasons. First of all, we want to focus on Vietnamese correction, instead of general text correction and learn to correct text based on context. Secondly, adding correction for English tremendously increase vocabulary size, which may leads to slower convergence time.

Result on Viwiki

Our model performance drop tremendously when tested against Viwiki dataset by Zalo. The research group crawled data from Wikipedia, a free to edit data source whose content is created by volunteers around the world, hence its contents are less monitored and checked for spelling. Nonetheless, our closer inspection reviews that this test set contains some interesting characteristics compared to VSEC.

Firstly, the research team gathered text from different various topics such as Vietnamese historical paragraph, chemistry, technical site building, etc. which contains many domain-specific terms. For example, chemical formulas such as "nitroglycerol" and "trinitrolycerin" or English words can easily confuses both the tokenizer and correction model. In historical texts, our model does not perform well for ancient phrases and human names. We see this as room for improvements moving forward to improve the quality of our training set and expand to various domain, as well as looking for ways to incorporate entity knowledge, such as names, organization, to our model.

Historical text	<p>Incorrect: Dụ Tông bèn xuống chiếu cho An phủ sứ các lộ đem quân các đội phong đoàn đi bắt giặc cướp, đến năm 1360 thì Bệ bị giết, đồng đảng hơn 30 người đều bị xử tử.</p> <p>Corrected: Dụ Tông bèn xuống chiếu cho An phủ sứ các lộ đem quân các đội phong đoàn đi bắt giặc cướp, đến năm 1360 thì Bệ bị giết, đồng đảng hơn 30 người đều bị xử tử.</p> <p>Predicted: Dụ Tông bèn xuống chiếu cho An phủ sứ các lộ đem quân các đội binh đoàn đi bắt giặc cướp, đến năm 1360 thì Bệ bị giết, đồng đảng hơn 30 người đều bị xử tử.</p>
Chemical names	<p>Incorrect: Các "khoáng vật cacbonat " bao gồm các khoáng vật chứa anion (CO₃)²⁻</p> <p>Corrected: Các "khoáng vật cacbonat " bao gồm các khoáng vật chứa anion (CO₃)²⁻</p> <p>Predicted: Các "khoáng vật cacbonat " bao gồm các khoáng vật chứa anion (CỎ)²⁻</p>

Table 7.1: Example of domain-specific keyword

Secondly, when looking the hand-labeled annotations, we notice that this test set contains error correction beyond our scope of Vietnamese error. For example, it also corrects English words and names such as "Harrrry" to "Harry". Also, it includes a vague shortcut correction, for example: "ko" is corrected to "không", while "ts.", which stands for "tiến sĩ", are not corrected. Another type of error is a insert token error.

English name	Incorrect: Kể từ khi đưa Harry vào thế giới pháp thuật, Hagrid đã trở thành một trong những người bạn thân thiết nhất của Harry. Corrected: Kể từ khi đưa Harry vào thế giới pháp thuật, Hagrid đã trở thành một trong những người bạn thân thiết nhất của Harry. Predicted: Kể từ khi đưa Harry vào thế giới pháp thuật, Hagrid đã trở thành một trong những người bạn thân thiết nhất của Harry.
Shortcuts	Incorrect: ... và bảo rằng cậu phải tìm cách trở về dòng thời gian Beta, nơi Mayuri sẽ ko chết. Corrected: ... và bảo rằng cậu phải tìm cách trở về dòng thời gian Beta, nơi Mayuri sẽ không chết. Predicted: ... và bảo rằng cậu phải tìm cách trở về dòng thời gian Beta, nơi Mayuri sẽ ko chết.

Table 7.2: Example of English word and Shortcut

Finally, we note that although the text were hand-labeled by separate teams, we still manage to find some diacritic and spelling error during the process. This affect our false positive rate, which in terms decrease model performance. Ultimately, this test set does not reflect our true target. Nevertheless, we plan to expand our model scope to other type of errors and tackle this test set in a different paper.

Initial approach for improving Viwiki test set score

On several attempt, we try to incorporate a pretrained named entity recognition (NER) system for increasing detection score. Our approach is to ignore person name, geographical location and organization. In particular, we use pretrained model provided by NLTK [47] and VnCoreNLP [48]. However, in testing, both pretrained model has several flaws while detecting entities.

- For NLTK package, we use the pretrained English model for detecting foreign language entities. Since the model is trained with English only, it often classify non-entities Vietnamese words as person names and organization.
- As for VnCoreNLP package, the model works incredibly well since it is trained on Vietnamese dataset. Though, in several cases where the names contains errors such as "Nguyenx Xuân Phúc", ignoring these errors would have affected our model overall performance. Another pretrained Vietnamese language model was intentionally used for scoring the sentence fluency whether we ignore or fix these errors would be beneficial. Nonetheless, publicly available Vietnamese pretrained language model is limited, we could not find any acceptable model that could be used. Not to mention, relying on too many models does not always yield good results.

Future works

Noticing our drawbacks, we have planned out several works that needed to be addressed for improving the performance.

- On the dataset level, expanding the size of dataset with more topics while incorporating more type of errors is a must to apprehend better real world data.
- Besides raw text features, integrating entities information or part of speech tags would be beneficial for model to learn more knowledge. Hence, enhancing the detection and correction performance.

- The problem of imbalanced classification remained on the detection phase, where most of the tokens are not erroneous. Applying more techniques is required as to improve precision.
- Researching on new features such as detecting sentence boundary, correct punctuation marks, and add missing words that match the sentence context where possible.

Chapter 8

Summary

In this paper, we tackle our goal for the thesis which is Vietnamese Contextual Spelling Correction using a deep-learning approach. We have demonstrated that leveraging a pre-trained transformer encoder such as BERT yields effective results for the task. We also enhance the correction model by incorporating a tokenization repair module, albeit accuracy is left to be improved. Our model scores well against the dataset released by Vietnam National University research group, outperforming paper's result in term of f-score. It, however, achieve poor score against public dataset by Zalo, which we mentioned contains errors beyond our scope.

For further research, we plan to study and extend our augmentation to capture more types of spelling errors or generalize our scope to textual correction. We would like to also explore options to incorporate domain knowledge to our model for inference.

References

- [1] Hieu Tran et al. *Hierarchical Transformer Encoders for Vietnamese Spelling Correction*. 2021. doi: 10.48550/ARXIV.2105.13578. url: <https://arxiv.org/abs/2105.13578>.
- [2] Sai Muralidhar Jayanthi, Danish Pruthi, and Graham Neubig. “NeuSpell: A Neural Spelling Correction Toolkit”. In: *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*. Online: Association for Computational Linguistics, Oct. 2020, pp. 158–164. doi: 10.18653/v1/2020.emnlp-demos.21. url: <https://www.aclweb.org/anthology/2020.emnlp-demos.21>.
- [3] Keisuke Sakaguchi et al. *Robust Word Recognition via semi-Character Recurrent Neural Network*. 2016. doi: 10.48550/ARXIV.1608.02214. url: <https://arxiv.org/abs/1608.02214>.
- [4] Erik Jones et al. *Robust Encodings: A Framework for Combating Adversarial Typos*. 2020. doi: 10.48550/ARXIV.2005.01229. url: <https://arxiv.org/abs/2005.01229>.
- [5] Shaohua Zhang et al. *Spelling Error Correction with Soft-Masked BERT*. 2020. arXiv: 2005.07421 [cs.CL].
- [6] Jinhua Zhu et al. *Incorporating BERT into Neural Machine Translation*. 2020. arXiv: 2002.06823 [cs.CL].
- [7] Phuong H. Nguyen et al. “Vietnamese spelling detection and correction using Bi-gram, Minimum Edit Distance, SoundEx algorithms with some additional heuristics”. In: *2008 IEEE International Conference on Research, Innovation and Vision for the Future in Computing and Communication Technologies*. 2008, pp. 96–102. doi: 10.1109/RIVF.2008.4586339.
- [8] Huong Nguyen et al. “Using Large N-gram for Vietnamese Spell Checking”. In: *Advances in Intelligent Systems and Computing* 326 (Jan. 2015), pp. 617–627. doi: 10.1007/978-3-319-11680-8_49.
- [9] Dinh-Truong Do et al. *VSEC: Transformer-based Model for Vietnamese Spelling Correction*. 2021. arXiv: 2111.00640 [cs.CL].
- [10] Kyunghyun Cho et al. *Learning Phrase Representations using RNN Encoder-Decoder for Statistical Machine Translation*. 2014. arXiv: 1406.1078 [cs.CL].
- [11] Sepp Hochreiter and Jürgen Schmidhuber. “Long Short-term Memory”. In: *Neural computation* 9 (Dec. 1997), pp. 1735–80. doi: 10.1162/neco.1997.9.8.1735.
- [12] Mike Schuster and Kuldeep Paliwal. “Bidirectional recurrent neural networks”. In: *Signal Processing, IEEE Transactions on* 45 (Dec. 1997), pp. 2673–2681. doi: 10.1109/78.650093.

- [13] Abdel-rahman Mohamed et al. “Deep bi-directional recurrent networks over spectral windows”. In: *2015 IEEE Workshop on Automatic Speech Recognition and Understanding (ASRU)*. 2015, pp. 78–83. doi: 10.1109/ASRU.2015.7404777.
- [14] Zhiyong Cui et al. *Deep Bidirectional and Unidirectional LSTM Recurrent Neural Network for Network-wide Traffic Speed Prediction*. 2019. arXiv: 1801.02143 [cs.LG].
- [15] Ilya Sutskever, Oriol Vinyals, and Quoc V. Le. *Sequence to Sequence Learning with Neural Networks*. 2014. arXiv: 1409.3215 [cs.CL].
- [16] Seiya Tokui et al. “Chainer: A Deep Learning Framework for Accelerating the Research Cycle”. In: *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. ACM. 2019, pp. 2002–2011.
- [17] Jacob Devlin et al. *BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding*. 2019. arXiv: 1810.04805 [cs.CL].
- [18] Rico Sennrich, Barry Haddow, and Alexandra Birch. “Neural Machine Translation of Rare Words with Subword Units”. In: *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Berlin, Germany: Association for Computational Linguistics, Aug. 2016, pp. 1715–1725. doi: 10.18653/v1/P16-1162. url: <https://aclanthology.org/P16-1162>.
- [19] Philip Gage. “A new algorithm for data compression”. In: *The C Users Journal archive* 12 (1994), pp. 23–38.
- [20] Ashish Vaswani et al. *Attention Is All You Need*. 2017. arXiv: 1706.03762 [cs.CL].
- [21] Kyunghyun Cho et al. “On the Properties of Neural Machine Translation: Encoder–Decoder Approaches”. In: *Proceedings of SSST-8, Eighth Workshop on Syntax, Semantics and Structure in Statistical Translation*. Doha, Qatar: Association for Computational Linguistics, Oct. 2014, pp. 103–111. doi: 10.3115/v1/W14-4012. url: <https://aclanthology.org/W14-4012>.
- [22] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. *Neural Machine Translation by Jointly Learning to Align and Translate*. 2016. arXiv: 1409.0473 [cs.CL].
- [23] Minh-Thang Luong, Hieu Pham, and Christopher D. Manning. *Effective Approaches to Attention-based Neural Machine Translation*. 2015. arXiv: 1508.04025 [cs.CL].
- [24] Jianpeng Cheng, Li Dong, and Mirella Lapata. *Long Short-Term Memory-Networks for Machine Reading*. 2016. arXiv: 1601.06733 [cs.CL].
- [25] Matthew E. Peters et al. “Deep contextualized word representations”. In: *CoRR abs/1802.05365* (2018).
- [26] Alec Radford et al. “Language Models are Unsupervised Multitask Learners”. In: (2019).
- [27] Jeremy Howard and Sebastian Ruder. “Fine-tuned Language Models for Text Classification”. In: *CoRR abs/1801.06146* (2018).
- [28] Yonghui Wu et al. “Google’s Neural Machine Translation System: Bridging the Gap between Human and Machine Translation”. In: *CoRR abs/1609.08144* (2016).
- [29] Yinhan Liu et al. “RoBERTa: A Robustly Optimized BERT Pretraining Approach”. In: *CoRR abs/1907.11692* (2019).
- [30] Zhilin Yang et al. “XLNet: Generalized Autoregressive Pretraining for Language Understanding”. In: *CoRR abs/1906.08237* (2019).

- [31] Zihang Dai et al. “Transformer-XL: Attentive Language Models Beyond a Fixed-Length Context”. In: *CoRR* abs/1901.02860 (2019).
- [32] Zhenzhong Lan et al. “ALBERT: A Lite BERT for Self-supervised Learning of Language Representations”. In: *CoRR* abs/1909.11942 (2019).
- [33] Kenji Imamura and Eiichiro Sumita. “Recycling a Pre-trained BERT Encoder for Neural Machine Translation”. In: *Proceedings of the 3rd Workshop on Neural Generation and Translation*. Hong Kong: Association for Computational Linguistics, Nov. 2019, pp. 23–31. doi: 10.18653/v1/D19-5603. url: <https://aclanthology.org/D19-5603>.
- [34] Masahiro Kaneko et al. *Encoder-Decoder Models Can Benefit from Pre-trained Masked Language Models in Grammatical Error Correction*. 2020. arXiv: 2005.00987 [cs.CL].
- [35] Dat Quoc Nguyen and Anh Tuan Nguyen. “PhoBERT: Pre-trained language models for Vietnamese”. In: *CoRR* abs/2003.00744 (2020). arXiv: 2003.00744. url: <https://arxiv.org/abs/2003.00744>.
- [36] Peter Norvig. *How to Write a Spelling Corrector*. url: <http://norvig.com/spell-correct.html>.
- [37] Hannah Bast, Matthias Hertel, and Mostafa M. Mohamed. “Tokenization Repair in the Presence of Spelling Errors”. In: *Proceedings of the 25th Conference on Computational Natural Language Learning*. Association for Computational Linguistics, 2021. doi: 10.18653/v1/2021.conll-1.22. url: <https://doi.org/10.18653/v1/2021.conll-1.22>.
- [38] Thomas Wolf et al. “HuggingFace’s Transformers: State-of-the-art Natural Language Processing”. In: *CoRR* abs/1910.03771 (2019). arXiv: 1910.03771. url: <http://arxiv.org/abs/1910.03771>.
- [39] Karen Kukich. *Techniques for automatically correcting words in text* Share on. 1992. doi: 146370.146380.
- [40] Cao Hong Nga et al. “Deep Learning Based Vietnamese Diacritics Restoration”. In: *2019 IEEE International Symposium on Multimedia (ISM)*. 2019, pp. 331–3313. doi: 10.1109/ISM46123.2019.00074.
- [41] Thai-Hoang Pham, Xuan-Khoai Pham, and Phuong Le-Hong. *On the Use of Machine Translation-Based Approaches for Vietnamese Diacritic Restoration*. 2017. arXiv: 1709.07104 [cs.CL].
- [42] Hung Bui. “Vietnamese Diacritics Restoration Using Deep Learning Approach”. In: Nov. 2018, pp. 347–351. doi: 10.1109/KSE.2018.8573427.
- [43] Myle Ott et al. “fairseq: A Fast, Extensible Toolkit for Sequence Modeling”. In: *Proceedings of NAACL-HLT 2019: Demonstrations*. 2019.
- [44] Alexis Conneau et al. “Unsupervised Cross-lingual Representation Learning at Scale”. In: *CoRR* abs/1911.02116 (2019).
- [45] Diederik P. Kingma and Jimmy Ba. *Adam: A Method for Stochastic Optimization*. cite arxiv:1412.6980Comment: Published as a conference paper at the 3rd International Conference for Learning Representations, San Diego, 2015. 2014. url: <http://arxiv.org/abs/1412.6980>.
- [46] Paulius Micikevicius et al. *Mixed Precision Training*. 2017. doi: 10.48550/ARXIV.1710.03740. url: <https://arxiv.org/abs/1710.03740>.

- [47] Steven Bird and Edward Loper. “NLTK: The Natural Language Toolkit”. In: *Proceedings of the ACL Interactive Poster and Demonstration Sessions*. Barcelona, Spain: Association for Computational Linguistics, July 2004, pp. 214–217. url: <https://aclanthology.org/P04-3031>.
- [48] Thanh Vu et al. “VnCoreNLP: A Vietnamese Natural Language Processing Toolkit”. In: *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Demonstrations*. New Orleans, Louisiana: Association for Computational Linguistics, June 2018, pp. 56–60. doi: 10.18653/v1/N18-5012. url: <https://aclanthology.org/N18-5012>.