

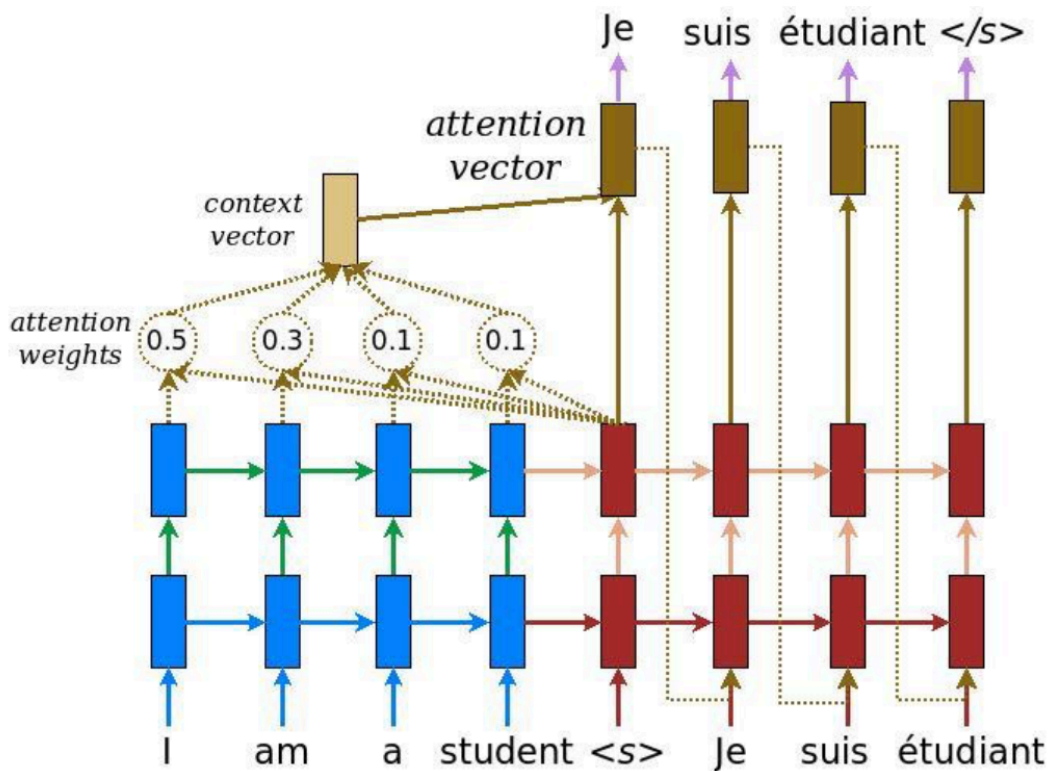
## CHƯƠNG 3: GIẢI PHÁP ĐỀ TÀI

### 3.1 TỔNG QUAN GIẢI PHÁP KIẾN TRÚC MÔ HÌNH

Nhóm sinh viên sử dụng mô hình kiến trúc đầu cuối (end-to-end hay còn được gọi với Sequence to Sequence Model – seq2seq) có ý tưởng từ bài báo *Sequence to Sequence Learning with Neural Networks* do nhóm tác giả đến từ google được ông bố vào năm 2014 tại Silicon Valley AI Lab đã trình bày ý tưởng cụ thể để xây dựng một mô hình mạng nơ-ron hồi quy tối ưu với hướng đi mới so với các hệ thống dịch máy truyền thống kết hợp cùng với cơ chế chú ý (Attention mechanism) từ bài báo *Effective Approaches to Attention-based Neural Machine* được thực hiện bởi nhóm tác giả đến từ đại học Stanford vào năm 2015. Từ đó nhóm sinh viên xây dựng một mô hình mạng nơ-ron hồi quy (Recurrent neural network) dùng để xây dựng một hệ thống dịch máy tiên tiến.

Mô hình seq2seq hoạt động dựa 2 mạng RNN kết hợp lại với một mạng RNN nhận nhiệm vụ mã hoá (encoder) câu đầu vào tiếng Anh thành một vector biểu diễn câu đầu vào và một mạng RNN giải mã (decoder ) có nhiệm vụ giải mã vector biểu diễn câu đầu vào và kết hợp với cơ chế chú ý (Attention mechanism) để giải mã câu đầu vào và cho ra kết quả câu đầu ra tiếng Việt tương ứng.

Kiến trúc tổng thể cho việc kết hợp 2 kiến trúc trên để xây dựng một mô hình dịch máy được minh hoạ cụ thể ở hình 3.1.



Hình 3.1: Tổng quan kiến trúc mô hình dịch máy

## 3.2 GIẢI PHÁP BIỂU DIỄN TỪ

### 3.2.1 Tổng quan về giải pháp

Để có thể sử dụng các mô hình Deep Learning (học sâu) phục vụ cho việc dịch máy, chúng ta cần biểu diễn các từ thành các số vì các mô hình chỉ làm việc với dữ liệu số. Vì thế dựa trên các kết quả tìm kiếm và thực nghiệm [], nhóm sinh viên đề xuất sử dụng Word Embedding (nhúng từ) dùng để biểu diễn các từ thành các vector số thực. Mô hình mà nhóm chọn là Word2vec với mục đích biểu diễn các từ tiếng Anh và tiếng Việt thành các vector số thực  $n$  chiều bằng nhau (mỗi chiều là một giá trị số thực) để phục vụ cho quá trình huấn luyện.

Word2vec là một mô hình học không giám sát (model unsupervised learning) nó dùng để thể hiện mối quan hệ giữa các từ, nó được kết hợp từ hai thuật toán Skip-gram và Continuous bag of words (CBOW). Ở đây nhóm sinh viên đề xuất sử dụng mô hình skip-gram cho biểu diễn từ. Với skip-gram, kích thước biểu diễn từ giảm từ kích thước bằng số từ trong bộ từ vựng xuống bằng chiều dài lớp ẩn. Hơn nữa các vector có ý nghĩa nhiều hơn về mặt mô tả mối quan hệ giữa các từ. Chi tiết mô hình đã được mô tả ở Chương 2.

### **3.2.2 Chi tiết giải pháp**

Nhóm sinh viên sử dụng đầu vào là tập dữ liệu được chia làm 2 tập tin chính chia làm 2 ngôn ngữ tiếng Anh và tiếng Việt.

Với bộ dữ liệu “IWSLT’15 English-Vietnamese data” với khoảng 100.000 câu song ngữ English-Vietnamese với dữ liệu thô chia làm hai tập tin tiếng English-Vietnamese như sau:

```
18 It 's a huge amount of stuff . It 's equal to the weight of methane .
19 And because it 's so much stuff , it 's really important for the atmospheric system .
20 Because it 's important to the atmospheric system , we go to all lengths to study this thing .
21 We blow it up and look at the pieces .
22 This is the EUPHORE Smog Chamber in Spain .
23 Atmospheric explosions , or full combustion , takes about 15,000 times longer than what happens in your car .
24 But still , we look at the pieces .
25 We run enormous models on supercomputers ; this is what I happen to do .
26 Our models have hundreds of thousands of grid boxes calculating hundreds of variables each , on minute timescales .
27 And it takes weeks to perform our integrations .
28 And we perform dozens of integrations in order to understand what 's happening .
29 We also fly all over the world looking for this thing .
30 I recently joined a field campaign in Malaysia . There are others .
```

```
18 Đó là một lượng khí thải khổng lồ , bằng tổng trọng lượng của mêtan .
19 Chính vì lượng khí thải rất lớn , nó có ý nghĩa quan trọng với hệ thống khí quyển .
20 Chính vì nó có ý nghĩa quan trọng với hệ thống khí quyển , giá nào chúng tôi cũng theo đuổi nghiên cứu này đến cùng .
21 Chúng tôi cho nó nổ và xem xét từng mảnh nhỏ .
22 Đây là Phòng nghiên cứu khói bụi EUPHORE ở Tây Ban Nha .
23 Nổ trong không khí hay cháy hoàn toàn diễn ra chậm hơn 15,000 lần so với những phản ứng trong động cơ xe .
24 Dù vậy , chúng tôi vẫn xem xét từng mảnh nhỏ .
25 Chúng tôi chạy những mô hình khổng lồ trên siêu máy tính ; đây là công việc của tôi .
26 Mô hình của chúng tôi gồm hàng trăm ngàn thùng xếp chồng tính toán với hàng trăm biến số trong thời gian cực ngắn .
27 Mà vẫn cần hàng tuần mới thực hiện xong các phép tích phân .
28 Chúng tôi cần làm hàng tá phép tính như thế để hiểu được những gì đang xảy ra .
29 Chúng tôi còn bay khắp thế giới để tìm phân tử này .
30 Gần đây tôi tham gia một cuộc khảo sát thực địa ở Malaysia . Còn nhiều chuyến khác nữa .
```

## ❖ Bước 1

Bước đầu tiên chúng ta thực hiện xử lí các câu dữ liệu như: xoá dấu “?”, “.”, xoá dấu khoảng trắng thừa và một số thứ khác.

Nhóm sinh viên thực hiện loại bỏ các câu có độ dài hơn 100.

## ❖ Bước 2

Nhóm sinh viên thực hiện tạo các từ điển word2int và int2word cho cả hai ngôn ngữ English-Vietnamese bằng cách sử dụng word\_tokenize và ta được kết quả như sau:

```
The word2index:
{'<pad>': 0, '<unk>': 1, '<s>': 2, '</s>': 3, '.': 4, ',': 5, 'tôi': 6, 'là': 7, 'và': 8, 'có': 9, 'một': 10}

The int2word:
{0: '<pad>', 1: '<unk>', 2: '<s>', 3: '</s>', 4: '.', 5: ',', 6: 'tôi', 7: 'là', 8: 'và', 9: 'có', 10: 'một'}
```

### ❖ Bước 3

Tiếp theo, nhóm sinh viên thực hiện chuyển từng câu song ngữ sang từng vector với từng từ ứng với vị trí của từ trong từ điển.

Với câu đầu vào tiếng Anh ta thực hiện thêm (padding) với những câu có độ dài bé hơn 100.

Với câu đầu vào tiếng Việt ta sẽ thực hiện thêm (padding) trong quá trình huấn luyện vì ta sẽ sử dụng độ dài thật của câu song ngữ để huấn luyện nhanh hơn. Và ta được kết quả như sau:

```
icle: the science behind a climate headline  
icle(int seq): [6, 310, 573, 13, 749, 4626, 0, 0, 0, 0, 0, 0, 0, 0,  
  
le: khoa học đằng sau một tiêu đề về khí hậu  
le(int seq): [326, 75, 1083, 116, 10, 372, 117, 41, 411, 743]
```

#### ❖ Bước 4

Tại đây ta thực hiện lấy nhúng từ (word embedding) tất cả các từ có trong từ điển English-Vietnamese.

[illegible]

### 3.3 GIẢI PHÁP XÂY DỰNG MÔ HÌNH DỊCH MÁÝ

Dựa trên các đánh giá thực tế và điều kiện phần cứng lẫn lượng dữ liệu (data) cho phép, nhóm sinh viên lựa chọn phương pháp học sâu (deep

learning) để xây dựng mô hình mạng nơ-ron hồi quy (Recurrent neural network) trong mô hình dịch máy (machine neural translation). Mô hình được đào tạo từ đầu đến cuối từ những câu đã được biểu diễn dưới các nhúng từ (word embedding) để tạo ra các chuỗi đầu vào bộ mã hoá (encoder) và bộ giải mã (decoder). Do đó với lượng dữ liệu đủ lớn và khả năng tính toán, mô hình có thể tự học một cách chính xác để thực hiện việc dịch một câu từ tiếng Anh sang tiếng Việt.

### 3.3.2 Mô hình mạng nơ-ron hồi quy (RNN) và khung huấn luyện

Cốt lõi của quá trình đào tạo một mô hình RNN là để nhận vào một văn bản tiếng Anh và tạo ra một văn bản tiếng Việt tương ứng. Để dễ hình dung ta có ví dụ một tập huấn luyện  $X = \{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots\}$  với  $x$  là một vector các nhúng từ (word embedding) tương ứng với câu tiếng Anh đầu vào và  $y$  là một nhãn tức là một vector các nhúng từ (word embedding) tương ứng với câu tiếng Việt ở đầu ra. Mỗi câu tiếng Anh  $x^{(i)}$  là một chuỗi thời gian có độ dài  $T^{(i)}$ , trong đó mỗi đoạn thời gian nhất định là một vector nhúng từ (word embedding)  $x_t^{(i)}$ ,  $t = 1, 2, \dots, T^{(i)}$ . Mục tiêu của RNN là chuyển đổi đầu vào  $x$  thành một chuỗi xác suất ký tự cho nhãn  $y$ , với  $y_t = P(w_t | x)$ , trong đó  $w_t$  thuộc các từ trong từ điển tiếng Việt và một vài ký tự đặc biệt khác.

Mô hình RNN được nhóm sinh viên chọn sử dụng là mô hình hồi quy với 3 thành phần chính. Thành phần đầu tiên là lớp nhúng từ (embedding), lớp (layer) này có nhiệm vụ chuyển các đầu vào của bộ mã hoá (encoder) và

bộ giải mã (decoder) từ dạng int sang dạng nhúng từ (word embedding) để phục vụ cho công việc tính toán phía sau.

Thành phần thứ hai là bộ mã hoá (encoder), với bộ mã hoá chúng ta sử dụng Multi layer Bi-directional LSTM với số lượng layer và số lượng hidden units của LSTM cell được thiết lập trong param. Ngoài ra nhóm sinh viên còn sử dụng DropoutWrapper để thiết lập giá trị Drop Out cho các LSTM cell để tránh hiện tượng quá khớp (over-fitting) với dữ liệu huấn luyện.

Thành phần thứ ba là bộ giải mã (decoder). Đối với bộ giải mã, chúng em chia thành hai trường hợp riêng biệt là huấn luyện mô hình(training) và dự đoán (inference). Trong quá trình huấn luyện chúng em sử dụng TrainingHelper còn khi dự đoán, chúng em sử dụng BasicDecoder với BeamSearchDecoder.

- ❖ Huấn luyện: chúng em sử dụng BahdanauAttention và TrainingHelper để huấn luyện mô hình. Chúng em còn sử dụng AdamOptimizer để cập nhật tham số cho mô hình và còn sử dụng Gradient Clipping để tránh mô hình bị bùng nổ độ dốc (exploding gradients).

- ❖ Dự đoán: sau khi huấn luyện xong mô hình và sử dụng mô hình này để dự đoán kết quả. Tuy nhiên do chúng ta không biết kết quả thực tế như trong quá trình huấn luyện, nên ta cần sử dụng các thuật toán tìm kiếm để cho ra kết quả phù hợp nhất và chúng em chọn sử dụng thuật toán tìm kiếm chùm tia (Beam Search) với beam-width = 10.



### **3.5 GIẢI PHÁP XÂY DỰNG MÁY CHỦ**

Máy chủ (server) được nhóm sinh viên chọn Amazon EC2 làm máy chủ với mục đích tạo ra một cầu nối giữa mô hình đã được huấn luyện (model) và phía ứng dụng sản phẩm (client) – được xây dựng với React Native. Vì vậy trong giới hạn của khoá luận, máy chủ chỉ cung cấp duy nhất một giao diện lập trình (API) với chức năng chuyển đổi từ một văn bản (text) tiếng Anh thành một văn bản (text) tiếng Việt tương ứng.

### **3.6 GIẢI PHÁP XÂY DỰNG ỨNG DỤNG**

Để ứng dụng hoá hệ thống dịch máy từ tiếng Anh sang tiếng Việt, nhóm sinh viên quyết định xây dựng web để ứng dụng kết quả của hệ thống vào một tình huống cụ thể có thể ứng dụng và thương mại hoá tốt.

Ứng dụng web do nhóm sinh viên xây dựng có chức năng chính là chuyển đổi văn bản tiếng Anh do người dùng nhập vào và đưa ra văn bản tiếng Việt tương ứng .

#### **3.6.1 Thiết kế giao diện ứng dụng**

Giao diện ứng dụng chỉ có một màn hình với chức năng chính là chuyển đổi một văn bản tiếng Anh thành một văn bản tiếng Việt tương ứng.

**Demo Mô Hình Dịch Máy Từ Tiếng Anh Sang Tiếng Việt**

| Tiếng Anh   | Tiếng Việt |
|---|------------|
| Nhập nội dung   |            |
| <div style="background-color: #007bff; color: white; padding: 5px 20px; display: inline-block; cursor: pointer;">Dịch</div> |            |

Hình 3.2 Màn hình chính của ứng dụng

Để sử dụng, người dùng nhập văn bản tiếng Anh vào ô tiếng Anh tương ứng và nhập vào nút dịch. Kết quả sẽ được hiển thị tại ô tiếng Việt.

### 3.6.2 Thiết kế kiến trúc ứng dụng

### 3.7 TỔNG KẾT

Thông qua chương 3, sinh viên đã làm rõ được các giải pháp cụ thể cho từng phần trong hệ thống dịch máy từ tiếng Anh sang tiếng Việt, hướng xây dựng máy chủ và cả ứng dụng trên nền tảng web.

Nhóm sinh viên đã trình bày một hệ thống dịch máy từ tiếng Anh sang tiếng Việt dựa trên việc học sâu (deep learning) từ đầu đến cuối có khả năng vượt trội và hiện đại trong hiện đại. Nhóm sinh viên tin rằng phương pháp này sẽ tiếp tục được cải thiện với các mô hình mới hơn, đơn giản hoặc phức tạp hơn khi tận dụng được sức mạnh tính toán phần cứng và kích thước dữ liệu được tăng thêm trong tương lai.

Chương kế tiếp nhóm sinh viên sẽ trình bày về các thư viện, công cụ và những khó khăn cụ thể nếu có cho các giải pháp đã trình bày ở chương này.