

Report of Neural Machine Translation Model

Huan Phan

VietAI Course 3's student, HCM, Vietnam

Abstract

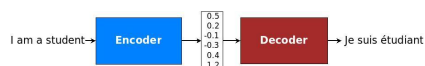
XYZ

Keywords: Neural Machine Translation, Seq2Seq, Attention

1. Overview of Neural Machine Translation

1.1. Encoder - Decoder Architect

Trong mô hình seq2seq dùng cho bài toán NMT (Neural Machine Translation) bao gồm 2 mạng RNN chính: Encoder và Decoder. Encoder với đầu vào là câu ở ngôn ngữ gốc, đầu ra tại layer cuối cùng của Encoder gọi là 1 context vector. Với ý nghĩa lượng thông tin từ câu của Encoder sẽ được tóm gọn lại trong 1 vector đầu ra cuối cùng. Từ đó, Decoder dùng chính context vector đó, cùng với hidden state và từ trước đó để predict từ tiếp theo tại decoder qua từng timestep.

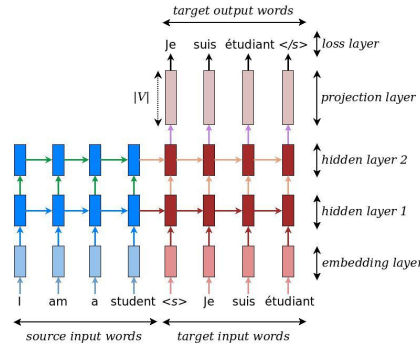


Hình 1: Encoder-decoder architecture – example of a general approach for NMT. An encoder converts a source sentence into a "meaning" vector which is passed through a decoder to produce a translation.

1.2. Neural Machine Translation Architect

1.3. Attention Mechanise

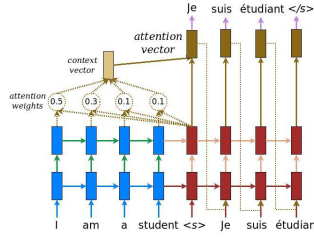
We conjecture that the use of a fixed-length vector is a bottleneck in improving the performance of this basic encoder–decoder architecture, and propose to extend this by allowing a model to automatically (soft-)search for parts of a source sentence that are relevant to predicting a target word, without having to form these parts as a hard segment explicitly[1]



Hình 2: Neural machine translation – example of a deep recurrent architecture proposed by for translating a source sentence "I am a student" into a target sentence "Je suis étudiant". Here, "<s>" marks the start of the decoding process while "</s>" tells the decoder to stop.

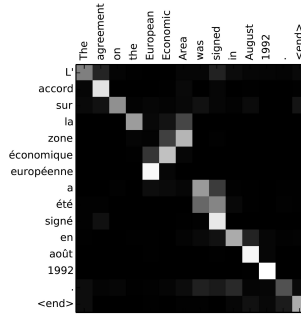
Việc encode toàn bộ thông tin từ source vào 1 vector cố định khiến việc mô hình khi thực hiện trên các câu dài (long sentence) không thực sự tốt, mặc dù sử dụng LSTM (BiLSTM, GRU) để khắc phục điểm yếu của mạng RNN truyền thống với hiện tượng Vanishing Gradient, nhưng như thế có vẻ vẫn chưa đủ, đặc biệt đối với những câu dài hơn những câu trong training data. Từ đó, trong paper, tác giả Bahdanau đề xuất 1 cơ chế cho phép mô hình có thể chú trọng vào những phần quan trọng (word liên kết với word từ source đến target), và thay vì chỉ sử dụng context layer được tạo ra từ layer cuối cùng của Encoder, tác giả sử dụng tất cả các output của từng cell qua từng timestep, kết hợp với hidden state của từng cell để "tổng hợp" ra 1 context vector (attention vector) và dùng nó làm đầu vào cho từng cell trong Decoder.

Cơ chế "tổng hợp" Attention trong paper của tác giả Bahdanau: *Align and Jointly model (Additive Attention)*.



Hình 3: Attention mechanism – example of an attention-based NMT system as described in (Luong et al., 2015) . We highlight in detail the first step of the attention computation. For clarity, we don't show the embedding and projection layers in Figure (2).

31 Ma trận bất đối xứng (confusion matrix) được tạo ra bởi alignment score,
32 thể hiện mức độ tương quan correlation giữa source và target.



Hình 4: Attention visualization – example of the alignments between source and target sentences. Image is taken from (Bahdanau et al., 2015).

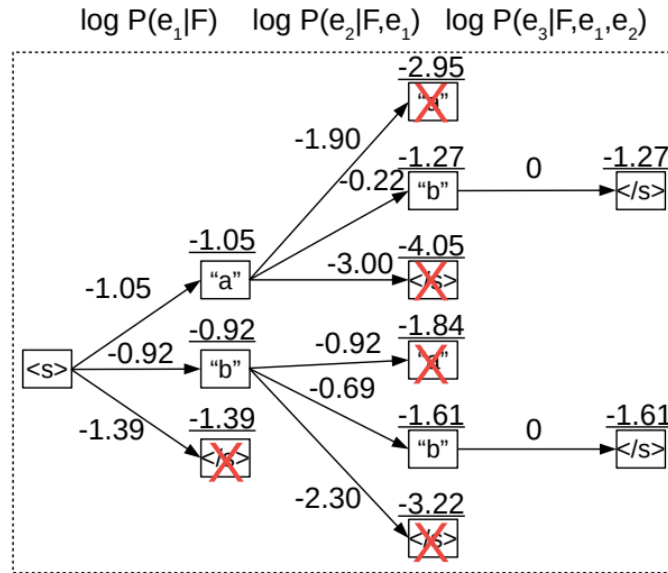
33 1.4. (SGD, Adam)

34 1.5. Inference mode

35 Được sử dụng trong quá trình decoding tại module decoder.

36 **Greedy:** Câu nguồn được module encoder encode thành encoder_state
37 để khởi tạo giá trị đầu vào cho module decoder. Quá trình decoding (dịch)
38 bắt đầu khi decoder đọc được ký tự bắt đầu câu $<s>$. Tại mỗi thời điểm t ,
39 output của mạng RNN là một vector V chiều (V là kích thước bộ vocab), với
40 mỗi phần tử của vector là các giá trị xác suất tương ứng với các từ trong
41 vocab. Chiến thuật Greedy (tham lam) là chọn từ tương ứng với vị trí của
42 xác suất cao nhất, từ này sẽ là output của NMT tại thời điểm t , và sẽ thành
43 đầu vào cho decoder tại thời điểm tiếp theo.

44 **Beam search:** Beam search decoder có thể làm tăng hiệu quả của mô
 45 hình. Ý tưởng của beam search là không chọn duy nhất một từ có giá trị xác
 46 suất dự đoán cao nhất (như greedy), mà giữ lại một tập hợp nhỏ top các từ
 47 có giá trị xác suất cao nhất. Tập các ứng cử viên này trở thành đầu vào cho
 48 decoder, tạo ra không gian output lớn hơn cho bước tiếp theo. Kích thước
 49 của tập hợp này gọi là beam width, thông thường có giá trị bằng 10, là đủ
 50 hiệu quả. Phương pháp này cho phép bước sau chọn ra ứng cử viên có xác
 51 suất đồng thời cao nhất từ các ứng cử viên đã chọn được từ bước trước, giảm
 52 thiểu lỗi chọn sai vì top các ứng cử viên từ bước trước đã được giữ lại.[2]



Hình 5: Ví dụ beam search với k=2. Trên các mũi tên là log probability của từng từ $P(e_t|F, e_1^{t-1})$. Trên các node là log probability tổng của các hypothesis tính đến node đó.

53 2. Measure performance with BLEU score

$$P(e_t|F, e_1^{t-1}) \quad (1)$$

54 3. Training with different hyper-parameters

55 3.1. IWSLT English-Vietnamese dataset

56 Train: 133K examples, vocab=vocab.(vi|en), train=train.(vi|en) dev=tst2012.(vi|en),
57 test=tst2013.(vi|en)

58 3.2. Data Input Pipeline

59 <https://github.com/tensorflow/nmtdata-input-pipeline>

60 3.3. Impact of choosing hyper-parameters

61 3.3.1. Training details

62 I train 2-layer LSTMs of 512 units with bidirectional encoder (i.e., 1
63 bidirectional layers for the encoder), embedding dim is 512. LuongAttention
64 (scale=True) is used together with dropout *keep_prob* of 0.8. All parameters
65 are uniformly. We use SGD with learning rate 1.0 as follows: train for 12K
66 steps (12 epochs); after 8K steps, we start halving learning rate every 1K
67 step.

68 Training system: Tesla K80

69 3.3.2. Result with VI - EN translation model

Parameter	Time	BLEU score
Treatment 1	0.0003262	0.562
Treatment 2	0.0015681	0.910
Treatment 3	0.0009271	0.296

Bảng 1: VI - EN performances

70 3.3.3. Result with EN - VI translation model

Experiment	Training Time	tst2012 (dev) BLEU	test2013 (test) BLEU
Treatment 1	0.0003262	0.562	1
Treatment 2	0.0015681	0.910	1
Treatment 3	0.0009271	0.296	1

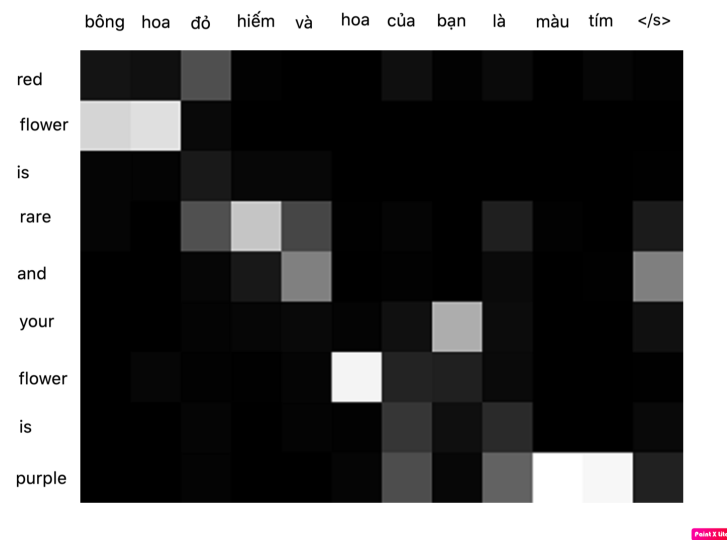
Bảng 2: VI - EN performances

71 3.3.4. Conclusion

72 4. Result visualization with Attention Matrix

73 4.1. Test 1

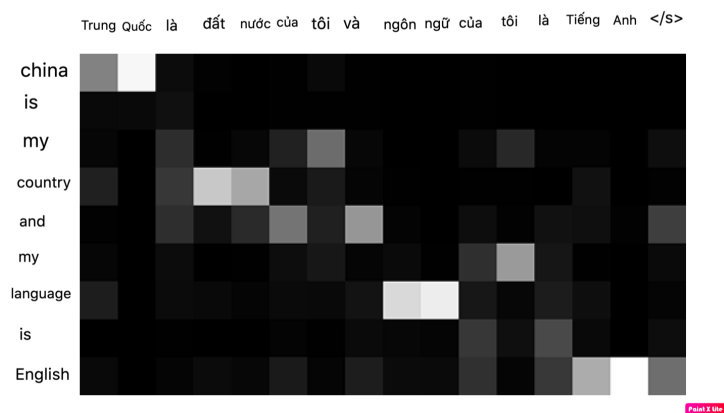
- 74 • red flower is rare and your flower is purple
- 75 • bông hoa đỏ hiếm và hoa của bạn là màu tím



Hình 6: Attention alignment giữa câu nguồn "red flower is rare and your flower is purple" và câu đích "bông hoa đỏ hiếm và hoa của bạn là màu tím".

76 4.2. Test 2

- 77 • China is my country and my language is English
- 78 • Trung Quốc là đất nước của tôi và ngôn ngữ của tôi là tiếng Anh



Hình 7: Attention alignment giữa câu nguồn "China is my country and my language is English" và câu đích "Trung Quốc là đất nước của tôi và ngôn ngữ của tôi là tiếng Anh".

79 Tài liệu

- 80 [1] J. M. Smith, A. B. Jones, Book Title, Publisher, 7th edition, 2012.
- 81 [2] Graham, Neural Machine Translation and sequence-to-sequence models:
- 82 A tutorial (????).