

Understanding and Improving Encoder Layer Fusion in Sequence-to-Sequence Learning



澳門大學
UNIVERSIDADE DE MACAU
UNIVERSITY OF MACAU

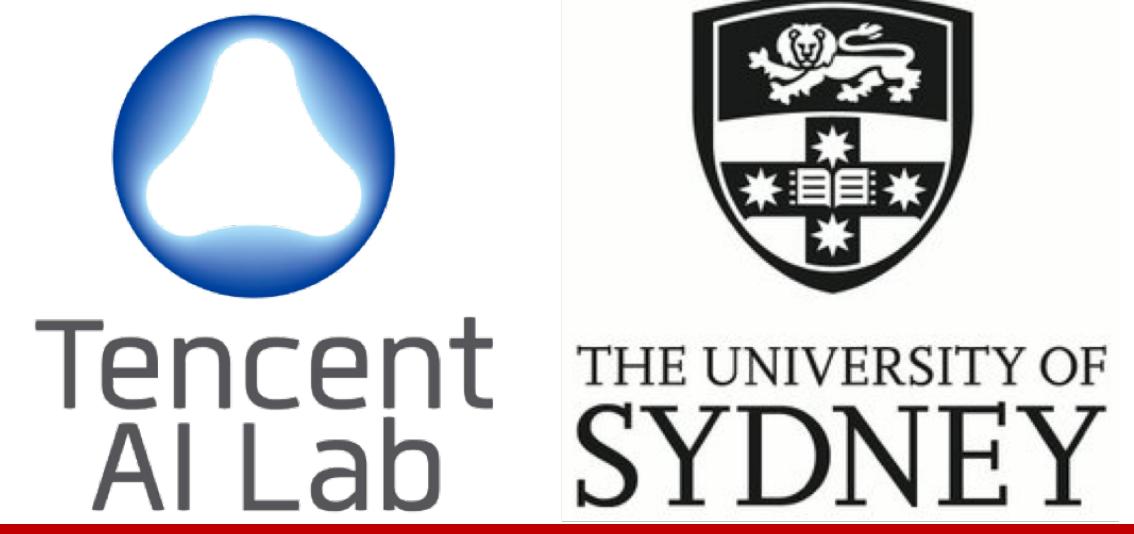


Xuebo Liu¹, Longyue Wang², Derek F. Wong¹, Liang Ding³, Lidia S. Chao¹, Zhaopeng Tu²

¹NLP²CT Lab, Department of Computer and Information Science, University of Macau

²Tencent AI Lab

³The University of Sydney



Introduction

Problems

- Recent studies reveal that fusing the intermediate encoder layers (EncoderFusion) is beneficial for sequence-to-sequence (Seq2Seq) models, such as layer aggregation [1], layer-wise coordination [2] and layer attention [3]. Despite its effectiveness, not much is known about how fusing encoder layer representations work.
- Other studies show that attending to lower encoder layers (excluding the encoder embedding layer) does not improve model performance [4], which is conflicted with existing conclusions. It is still unclear why and when fusing encoder layers should work in Seq2Seq models.

Contributions

- Introduce a fine-grained layer attention method to qualitatively and quantitatively evaluate the contribution of individual encoder layers.
- Demonstrate that the encoder embedding layer is essential for fusing encoder layers, which consolidates conflicted findings reported by previous studies.
- Propose a simple yet effective SurfaceFusion approach to directly exploit the encoder embedding layer for the decoder, which produces more expressive bilingual embeddings.

Source Representation Bottleneck

Seq2Seq decoder only takes the abstract representations at uppermost layer as input, which ignores other usefully surface representations at other layers. We hypothesize that its limited representation capacity may not sufficiently model those surface features from lower encoder layers, especially the embedding layer. We call such an issue as source representation bottleneck.

Fine-grained Layer Attention (FGLA)

We propose FGLA to investigate the source representation bottleneck. Specifically, FGLA replaces the layer-agnostic source representation X^N with the layer-aware representation S^m for each decoder layer Y^m :

$$S^m = \sum_{n=0}^N \hat{w}^{m,n} \odot X^n, \quad \hat{w}^{m,n} = [\hat{w}^{m,n,1}, \dots, \hat{w}^{m,n,D}], \quad \hat{w}^{m,n,d} = \frac{\exp(w^{m,n,d})}{\sum_{n'=0}^N \exp(w^{m,n',d})}$$

where \odot denotes an element-wise multiplication, and $w^{m,n,d}$ denotes an element in the learnable attention weight $W \in R^{M \times (N+1) \times D}$, where D is the dimensionality of the source representation.

Experimental Setup

Machine translation: WMT16 Romanian-English; WMT14 English-{German, French}

Text summarization: CNN/Daily Mail

Grammatical error correction: CONLL14

Understanding Encoder Layer Fusion (EncoderFusion)

Visualization of layer attention

In all tasks, higher decoder layers especially the uppermost ones pay more attention to the encoder embedding layer, which indicates that the surface representations potentially bring some additional useful features to the model performance.

Contribution of individual encoder layer

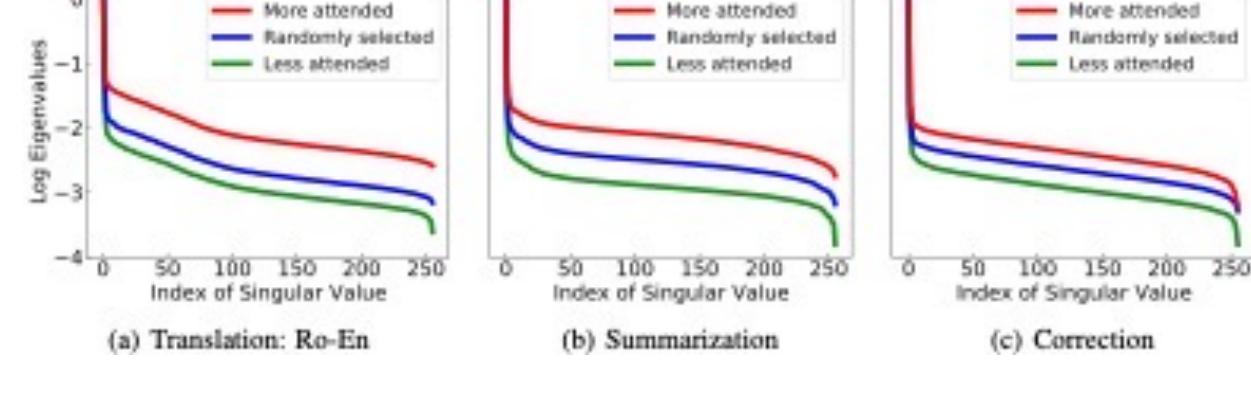
- Masking the encoder embedding layer seriously harms the model performance in all tasks.
- Masking the encoder embedding layer consistently increases the length of generated output.

Encoder Layer	Decoder Layer 1	Decoder Layer 2	Decoder Layer 3	Decoder Layer 4	Decoder Layer 5	Decoder Layer 6
Emb.	1.15, 1.18, 1.17, 1.16, 1.15, 1.14	5.12, 5.16, 5.18, 5.17, 5.15, 5.14	4.12, 4.14, 4.15, 4.15, 4.13, 4.11	3.12, 3.14, 3.13, 3.14, 3.11, 3.11	2.14, 2.12, 2.13, 2.14, 2.15, 2.12	1.16, 1.18, 1.19, 1.17, 1.16, 1.14
1	1.14, 1.11, 1.13, 1.15, 1.11, 1.12	4.14, 4.15, 4.15, 4.13, 4.10, 4.10	3.13, 3.14, 3.13, 3.11, 3.12, 3.11	2.14, 2.12, 2.13, 2.15, 2.16, 2.16	1.16, 1.13, 1.14, 1.17, 1.18, 1.19	1.13, 1.11, 1.10, 1.14, 1.18, 1.19
2	1.14, 1.11, 1.13, 1.15, 1.11, 1.12	4.14, 4.15, 4.15, 4.13, 4.10, 4.10	3.13, 3.14, 3.13, 3.11, 3.12, 3.11	2.14, 2.12, 2.13, 2.15, 2.16, 2.16	1.16, 1.13, 1.14, 1.17, 1.18, 1.19	1.13, 1.11, 1.10, 1.14, 1.18, 1.19
3	1.14, 1.11, 1.13, 1.15, 1.11, 1.12	4.14, 4.15, 4.15, 4.13, 4.10, 4.10	3.13, 3.14, 3.13, 3.11, 3.12, 3.11	2.14, 2.12, 2.13, 2.15, 2.16, 2.16	1.16, 1.13, 1.14, 1.17, 1.18, 1.19	1.13, 1.11, 1.10, 1.14, 1.18, 1.19

(a) Translation: Ro-En (b) Summarization (c) Correction



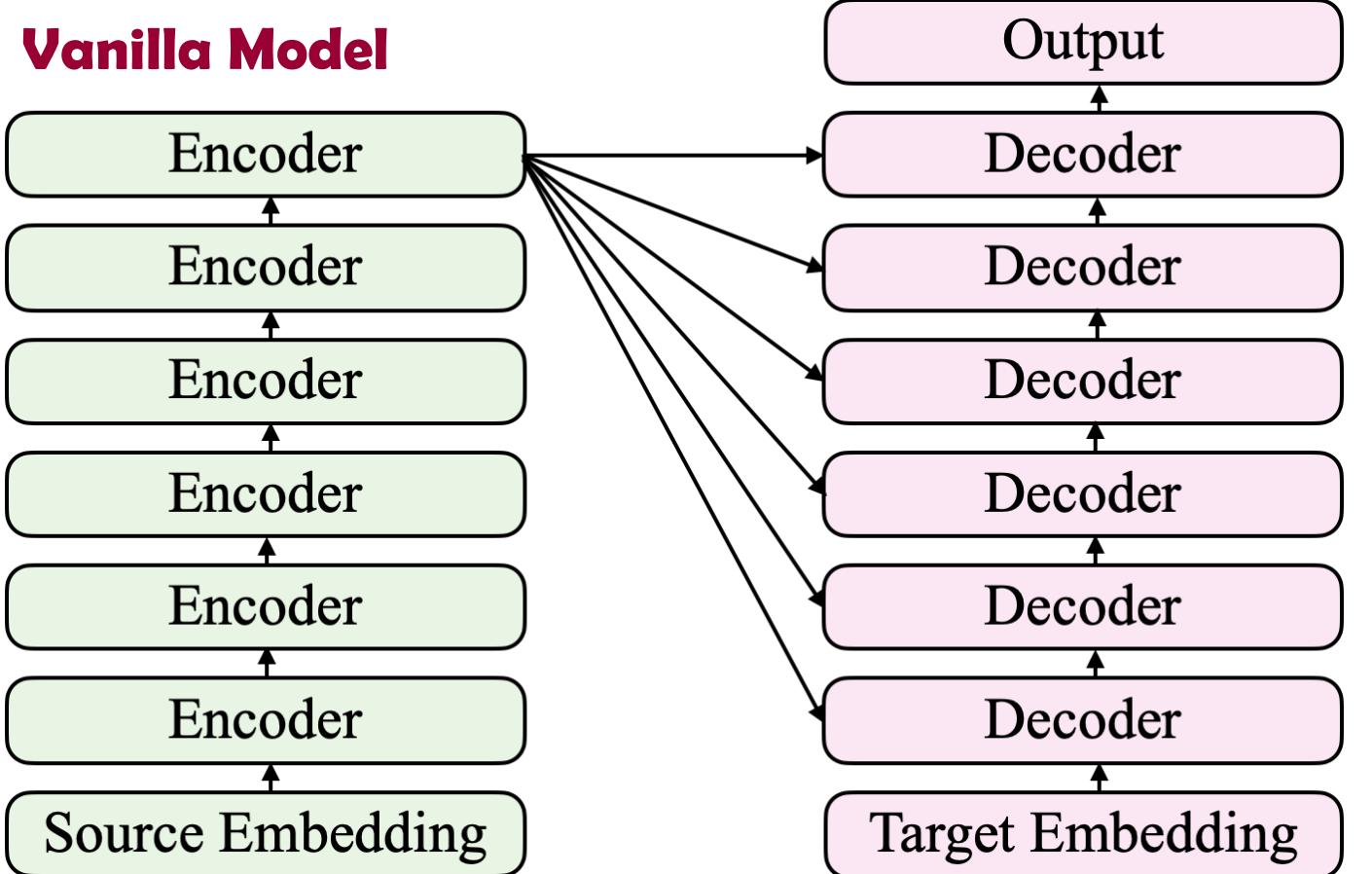
(a) Model performance (b) Length of output



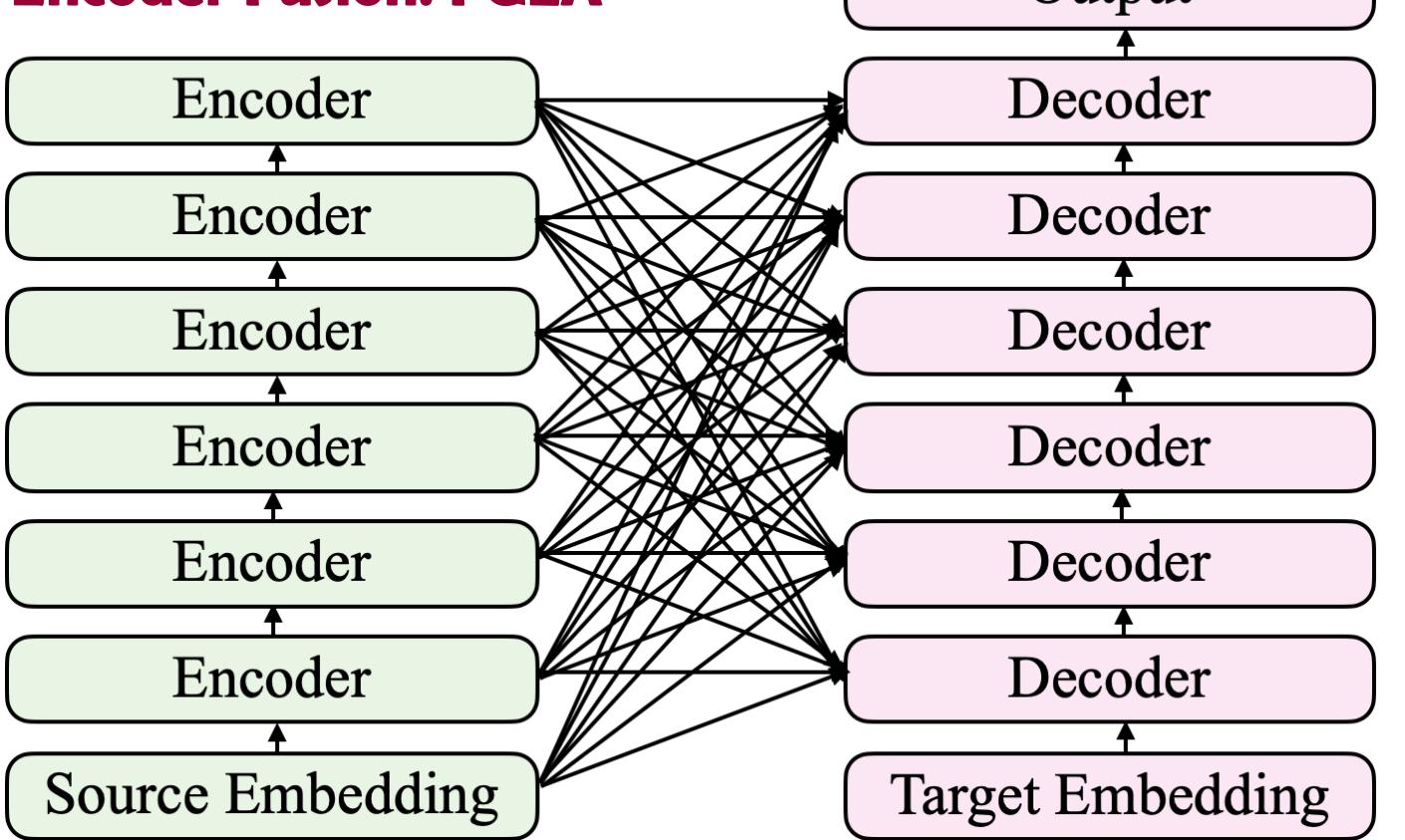
(a) Translation: Ro-En (b) Summarization (c) Correction

Vanilla Model vs. EncoderFusion

Vanilla Model



Encoder Fusion: FGLA



Improving EncoderFusion: SurfaceFusion

SurfaceFusion: Directly Exploit Source Embedding for Decoder

We rewrite $P(y_j)$ as a fused probability with the second condition term x :

$$\hat{y} = \arg \max \prod_{j=1}^J \log P(y_j) = \arg \max \prod_{j=1}^J \log \Phi(P(y_j|y_{<j}, x), P(y_j|x)) \\ P(y_j|x) = \frac{\exp(\mathbb{1}_{y_j}(\tilde{r}(y_j, x))/\tau)}{\sum_{w \in \mathcal{V}_y} \exp(\mathbb{1}_w(\tilde{r}(y_j, x))/\tau)}$$

where $\Phi(\cdot)$ is a fusion method. $\tilde{r}(y_j, x)$ denotes the surface feature calculated by an extra attention module and $\tilde{r}(y_j, x)$ denotes its logit version transformed by the pre-softmax weight. We use τ to control the sharpness of the probability distribution. As τ approaches to 0, the distribution tends to be a one-hot distribution representing the token of the maximum probability. The distribution becomes uniform at a higher τ .

Choices of Fusion Function

Hard Fusion: $\Phi_{\text{hard}} = \lambda \log P(y_j|y_{<j}, x) + (1 - \lambda) \log P(y_j|x)$

Soft Fusion: $\Phi_{\text{soft}} = \log(\text{softmax}(E(y_j|y_{<j}, x) + \log P(y_j|x))$

where λ is an interpolation weight, and $E(y_j|y_{<j}, x)$ is the pre-softmax logit of the probability $P(y_j|y_{<j}, x)$. Compared to hard fusion, soft fusion removes the need for manually setting the hyperparameter λ .

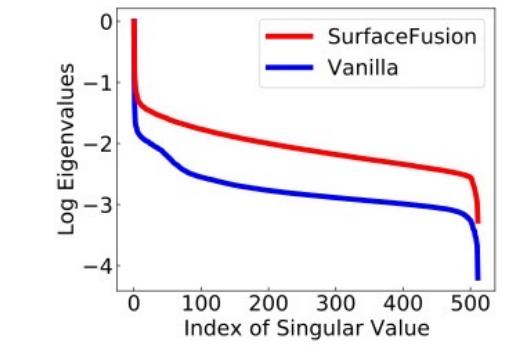
Closer Relationship between Word Embeddings

SurfaceFusion shortens the path distance between source and target embeddings, which can help to learn better bilingual embeddings with direct interactions. This can help the source and target embeddings to learn better representations, and has been validated to be beneficial for Seq2Seq models [5].

All	Non-Shared
Vanilla	0.602
SurfaceFusion	0.650

More Expressive Word Embeddings

The singular values of SurfaceFusion become more uniformly distributed, which demonstrates that the fused surface features enhance the representation learning of embeddings. This provides a better starting point for the model to effectively extract surface and abstract features, which leads to an improvement of model performance.



Main Results

	Translation			Summarization			Correction		
	Ro-En	En-De	En-Fr	RG-1	RG-2	RG-L	Prec.	Recall	F _{0.5}
Existing	34.0	29.3	43.2	40.1	17.6	36.8	65.5	33.1	54.8
Vanilla	33.8	28.9	43.4	40.4	17.7	37.2	64.7	33.2	54.3
FGLA	34.5	29.1	43.5	40.8	18.0	37.5	67.7	31.9	55.3
Hard fusion	35.1	29.5	43.9	40.9	18.2	37.7	67.0	34.4	56.3
Soft fusion	34.0	29.0	43.6	41.0	18.3	37.9	66.8	35.0	56.6

Acknowledgement

This work was supported in part by the Science and Technology Development Fund, Macau SAR (Grant No. 0101/2019/A2), and the Multi-year Research Grant from the University of Macau (Grant No. MYRG2020-00054-FST). We thank the anonymous reviewers for their insightful comments.

References

- Zi-Yi Dou, Zhaopeng Tu, Xing Wang, Shuming Shi, and Tong Zhang. Exploiting deep representations for neural machine translation. In EMNLP, 2018.
- Tianyu He, Xu Tan, Yingce Xia, Di He, Tao Qin, Zhibo Chen, and Tie-Yan Liu. Layer-wise coordination between encoder and decoder for neural machine translation. In NIPS, 2018.
- Ankur Bapna, Mia Xu Chen, Orhan Firat, Yuan Cao, and Yonghui Wu. Training deeper neural machine translation models with transparent attention. In EMNLP, 2018.
- Tobias Domhan. How much attention do you need? a granular analysis of neural machine translation architectures. In ACL, 2018.
- Xuebo Liu, Derek F. Wong, Yang Liu, Lidia S. Chao, Tong Xiao, and Jingbo Zhu. Shared-private bilingual word embeddings for neural machine translation. In ACL, 2019.