

# Von Mises–Fisher Loss For Training Sequence To Sequence Models With Continuous Outputs

ICLR 2019

Kunmar et al.

---

Presenter: Jungsoo Park

Data Mining & Information Systems Lab.  
Department of Computer Science and Engineering,  
College of Informatics, Korea University

---

## Time Matters

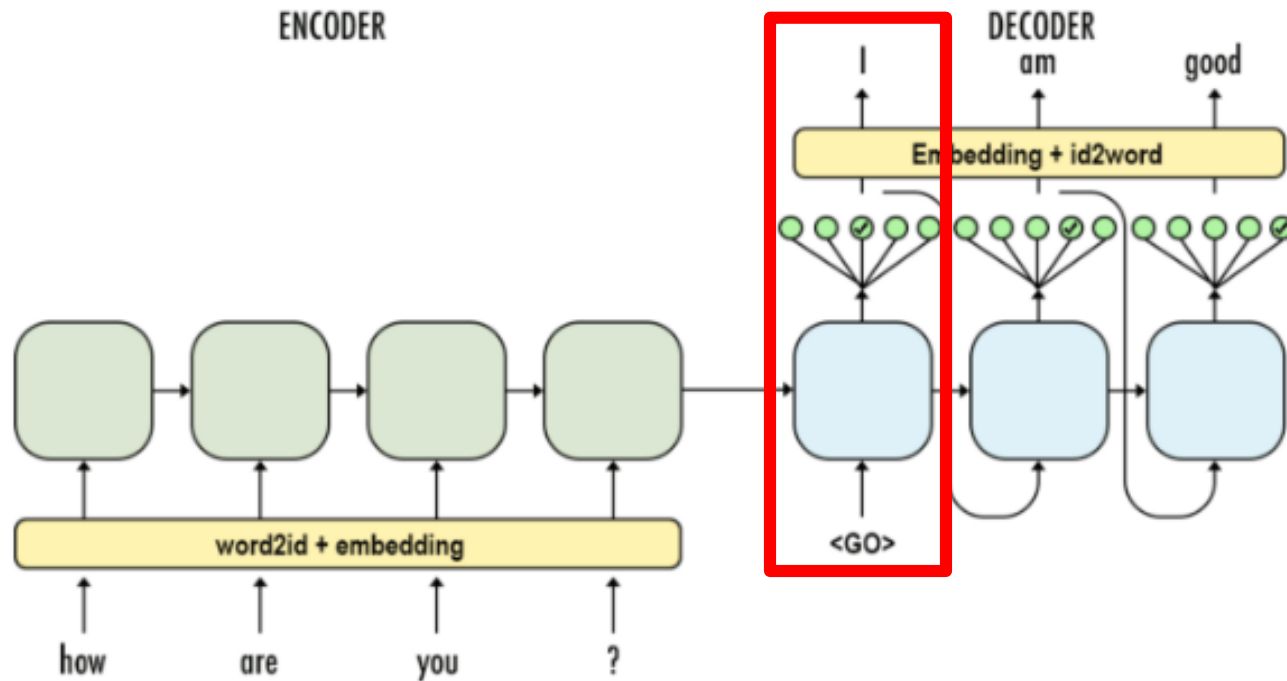


As for BERT  
Pre-Training took 4 days using

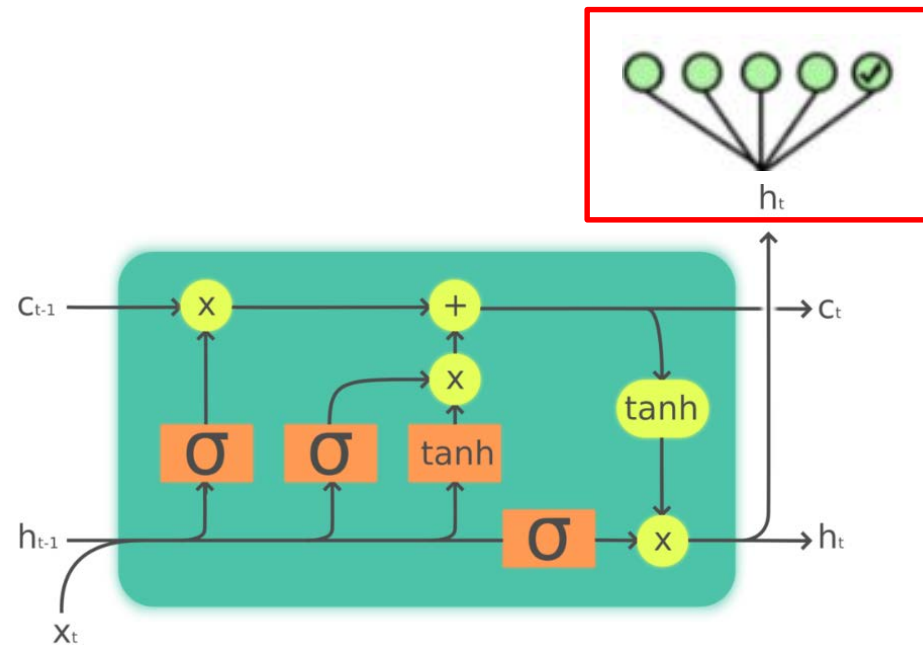
- BERT-Base: 4 Cloud TPUs (16 TPU chips total)
- BERT-Large: 16 Cloud TPUs (64 TPU chips total)

If using 8 TESLA P100,  
it would take **one year in training**

## Seq2Seq



## Decoding



Bottleneck in training

$$s_w = W_{hw} \mathbf{h}_t + b_w$$

$$W \in \mathbb{R}^{V \times H}$$

$$b \in \mathbb{R}^v$$

Bottleneck in training

$$\mathbf{p}_t(w) = \frac{e^{s_w}}{\sum_{v \in \mathcal{V}} e^{s_v}}$$

$$\text{NLL}(\mathbf{p}_t, \mathbf{o}(w)) = -\log(\mathbf{p}_t(w))$$

## Byte Pair Encoding

Initial Vocab    `low : 5, lower : 2, newest : 6, widest : 3`

Characterize

```
# dictionary
low : 5, lower : 2, newest : 6, widest : 3
```

```
# vocabulary
l, o, w, e, r, n, w, s, t, i, d
```

```
low : 5,
lower : 2,
newest : 6,
widest : 3
```



```
low : 5,
lower : 2,
newest : 6,
widest : 3
```



```
low : 5,
lower : 2,
newest : 6,
widest : 3
```

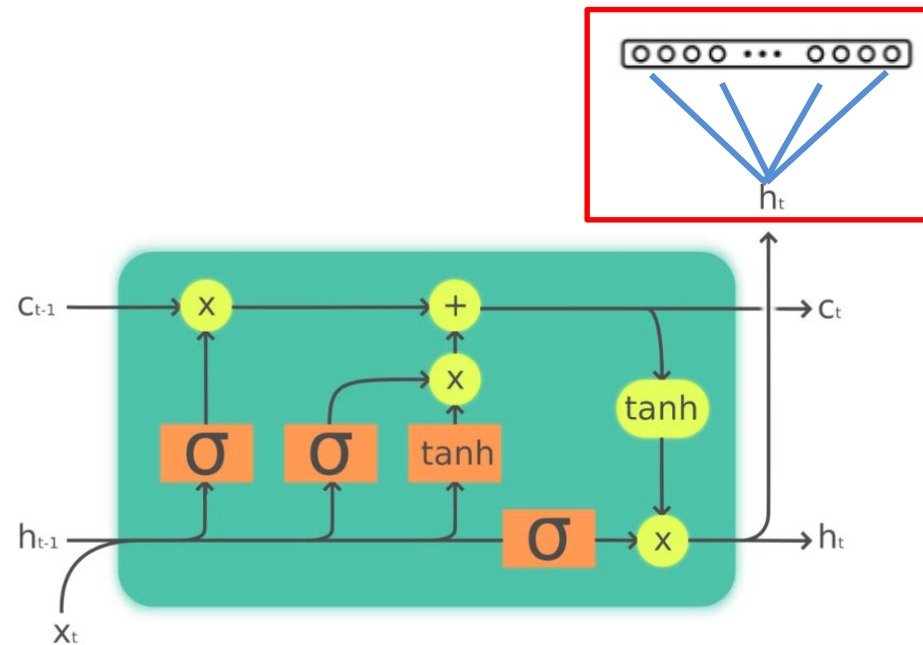


```
low : 5,
lower : 2,
newest : 6,
widest : 3
```

**# vocabulary update!**

l, o, w, e, r, n, w, s, t, i, d, es, est lo, low, ne, new, newest, wi, wid, widest

## Decoding

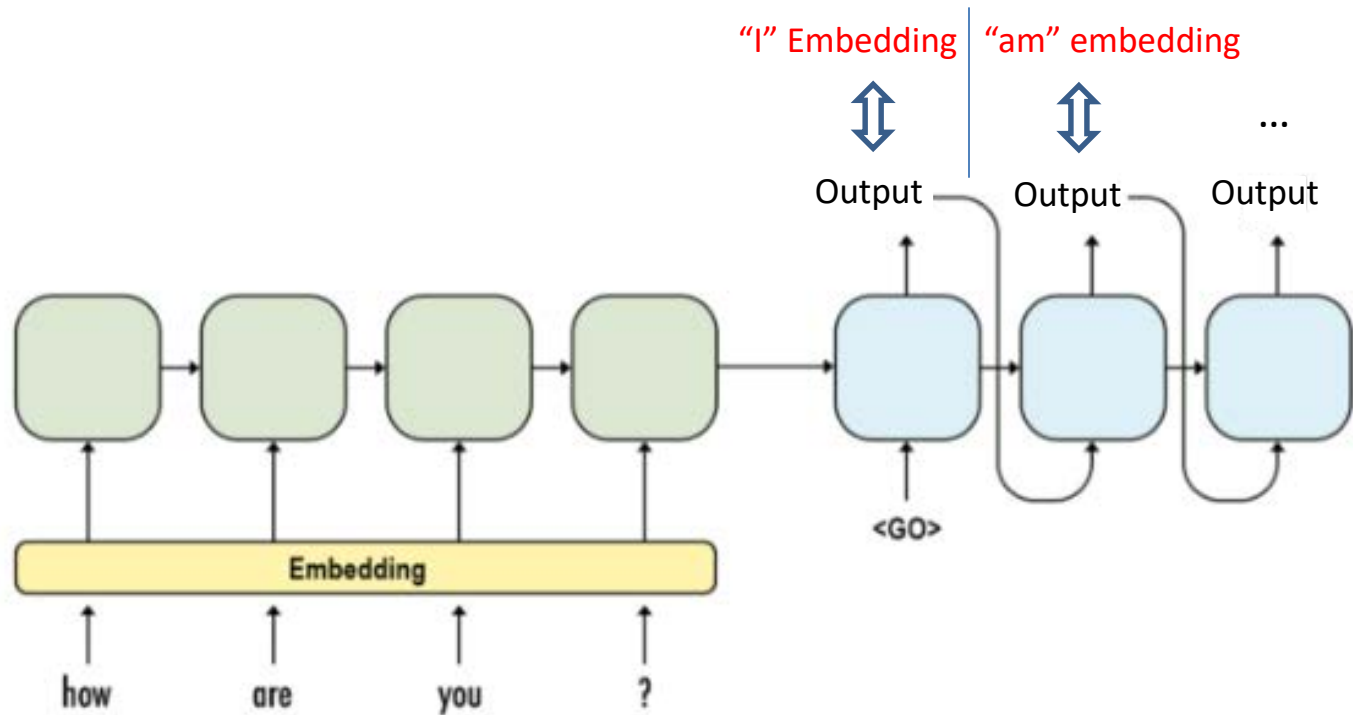


Predicting  
Embedding Vector

No need for soft-max operation, but  
predict the word embedding vector itself

## Overview

Loss occurs from the **target word(label's)**  
**embedding vector** and our generated one's



## Pros

1. Big Batch Size
2. Unique Approach to Handle Bottleneck in Softmax

## Cons

1. Beam Search Not Available (KNN)
2. Inference Not Efficient
3. Lack of Experiments based on Contextualized Representation

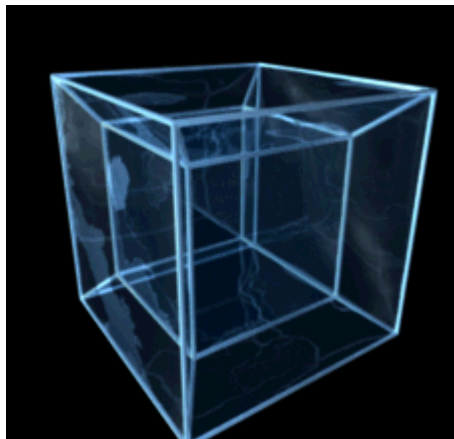


## Loss Function

But what **loss function** should we use?

**L2 Loss?**

Basic assumption behind using L2 loss is that  
the output space follows **Gaussian distribution**  
and is a **Distance-based metric**



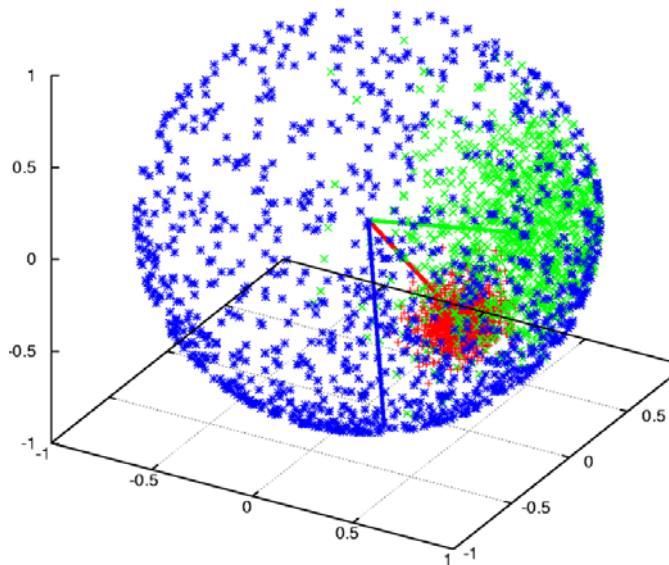
Hypercube

“Most of the volume of a  
high-dimensional orange is  
**in the skin, not the pulp.**”

Therefore, it's not the best way to use distance metric,  
but instead use **direction based metric** for high dimension data

### Von Mises-Fisher

Directional Equivalent of Gaussian  
Defined on the Support of (p-1)  
Dimensional Sphere



$$p(\mathbf{e}(w); \boldsymbol{\mu}, \kappa) = C_m(\kappa) e^{\kappa \boldsymbol{\mu}^T \mathbf{e}(w)},$$

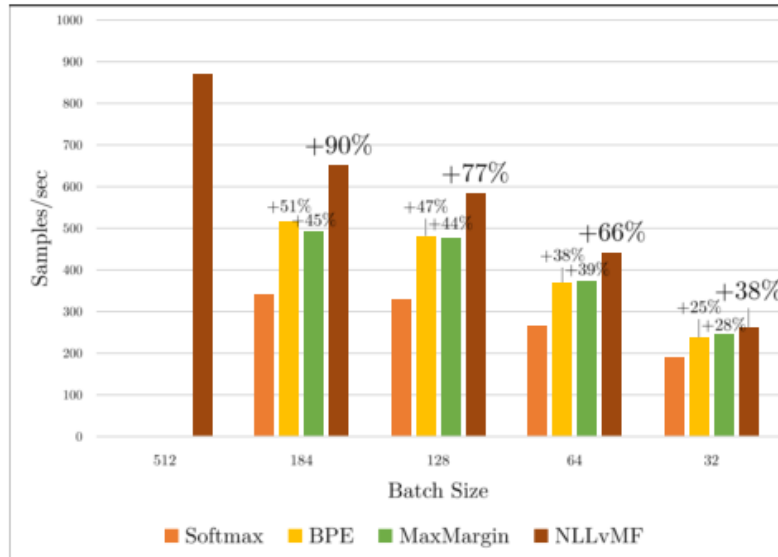
$$C_m(\kappa) = \frac{\kappa^{m/2-1}}{(2\pi)^{m/2} I_{m/2-1}(\kappa)},$$

$$p(\mathbf{e}(w); \hat{\mathbf{e}}) = \text{vMF}(\mathbf{e}(w); \hat{\mathbf{e}}) = C_m(\|\hat{\mathbf{e}}\|) e^{\hat{\mathbf{e}}^T \mathbf{e}(w)}$$

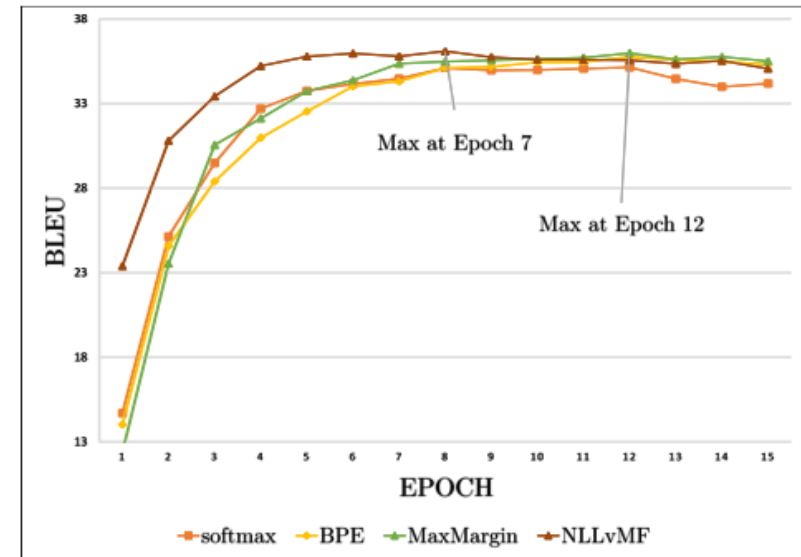
$$\text{NLLvMF}(\hat{\mathbf{e}}; \mathbf{e}(w)) = -\log(C_m(\|\hat{\mathbf{e}}\|)) - \hat{\mathbf{e}}^T \mathbf{e}(w)$$

Embedding Model	Tied Emb	Source Type/ Target Type	Loss	BLEU		
				fr-en	de-en	en-fr
-	no	word→word	CE	31.0	24.7	29.3
-	no	word→BPE	CE	29.1	24.1	29.8
-	no	BPE→BPE	CE	<b>31.4</b>	<b>25.8</b>	<b>31.0</b>
word2vec	no	word→emb	L2	27.2	19.4	26.4
word2vec	no	word→emb	Cosine	29.1	21.9	26.6
word2vec	no	word→emb	MaxMargin	29.6	21.4	26.7
fasttext	no	word→emb	MaxMargin	31.0	25.0	29.0
fasttext	yes	word→emb	MaxMargin	<b>32.1</b>	<b>25.0</b>	<b>31.0</b>
word2vec	no	word→emb	NLLvMF <sub>reg1</sub>	29.5	22.7	26.6
word2vec	no	word→emb	NLLvMF <sub>reg1+reg2</sub>	29.7	21.6	26.7
word2vec	yes	word→emb	NLLvMF <sub>reg1+reg2</sub>	29.7	22.2	27.5
fasttext	no	word→emb	NLLvMF <sub>reg1+reg2</sub>	30.4	23.4	27.6
fasttext	yes	word→emb	NLLvMF <sub>reg1+reg2</sub>	<b>32.1</b>	<b>25.1</b>	<b>31.7</b>

- x2.5 times faster than training word -> word (baseline) but shows comparable result
- Proposed Loss results in better result than the empirical loss functions(L2, Cosine Loss, Max Margin Loss)



Samples processed per second  
(Proposed method outperforms)



Convergence  
(NLLvMF Loss outperforms in stability)