

INTRODUCTION & OBJECTIVES

This study was done on translation task, and aimed to increase its performance by diversifying multi-head attention distributions and encouraging attention distributions to be sparse by incorporating Determinantal Point Processes.

- Objective 1) To generate sparse & diversified multi-head attention**
Objective 2) To enhance the performance of Neural Machine Translation

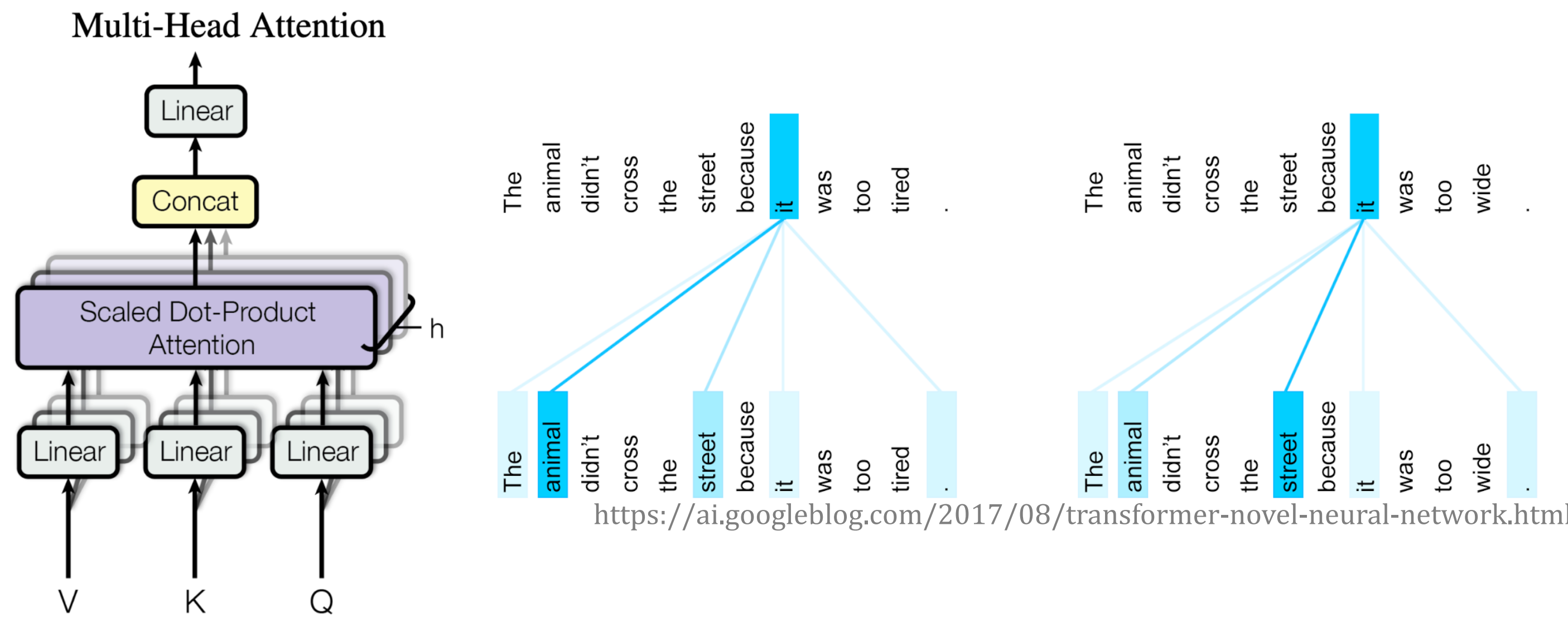
RELATED WORKS

Transformer

Attention is all you need. Vaswani et al. NeurIPS 2017

Multi-Head attention attends on different parts of sentences

- $Q^h, K^h, V^h = QW_h^Q, KW_h^K, VW_h^V$
- $O^h = A^h V^h$ with $A^h = \text{Attention}^h = \text{softmax}\left(\frac{Q^h K^{hT}}{\sqrt{d_k}}\right)$
- H Attention functions to produce $\{O^1, \dots, O^H\}$
- $\text{Multihead} = \text{Concat}(O^1, \dots, O^H)W^O$ (H : number of Attention head)



Transformer w/ disagreement regularization

Multi-Head Attention with Disagreement Regularization. Li et al. ENMLP 2018

Use Cosine-Similarity Between Distributions

- $J(\theta) = \text{argmax}_{\theta} \{L(y|x; \theta) + \lambda * D(a|x, y; \theta)\}$
- $D_{\text{subspace}} = -\frac{1}{H^2} \sum_{i=1}^H \sum_{j=1}^H \frac{v^i \cdot v^j}{\|v^i\| \|v^j\|}$
- $D_{\text{position}} = -\frac{1}{H^2} \sum_{i=1}^H \sum_{j=1}^H |A^i \odot A^j|$
- $D_{\text{output}} = -\frac{1}{H^2} \sum_{i=1}^H \sum_{j=1}^H \frac{o^i \cdot o^j}{\|o^i\| \|o^j\|}$

However, this baseline model has no guarantee to generate sparse attention distributions to attend well to specific words.

PROPOSED METHOD

Our proposed DPP based Multi-Head Attention:

1) Quality term of L-Matrix

- $q_i = \frac{1}{1 + \text{entropy}(A^i)}$: Generate sparse attention distributions

2) Diversity term of L-Matrix

- $D_{\text{subspace}} : \phi_i^T \phi_j = \frac{1}{H^2} \sum_{i=1}^H \sum_{j=1}^H \frac{v^i \cdot v^j}{\|v^i\| \|v^j\|}$
- $D_{\text{position}} : \phi_i^T \phi_j = \frac{1}{H^2} \sum_{i=1}^H \sum_{j=1}^H \frac{o^i \cdot o^j}{\|o^i\| \|o^j\|}$
- $D_{\text{output}} : \phi_i^T \phi_j = \frac{1}{H^2} \sum_{i=1}^H \sum_{j=1}^H \frac{A^i \cdot A^j}{\|A^i\| \|A^j\|}$

3) Constructing L-matrix

- Quality-Diversity Decomposition
- $L_{ij} = q_i \phi_i^T \phi_j q_j$ (q_i : quality score of item i)
($\phi_i^T \phi_j \in [-1, 1]$: normalized similarity between items i, j)

4) Loss Function

- $J(\theta) = \text{argmax}_{\theta} \{L(y|x; \theta) + \lambda * D_{\text{dpp}}\}$
- $D_{\text{dpp}} = \det(L)$ where $L_{i,j} = q_i \phi_i^T \phi_j q_j$

Determinantal Point Process

Determinantal point processes for machine learning. Kulesza et al.

Probability Measure P : $P_L(Y = Y) = \frac{\det(L_Y)}{\det(L+I)}$

1) Algebraic Intuition of DPP

➤ Quality of the Item

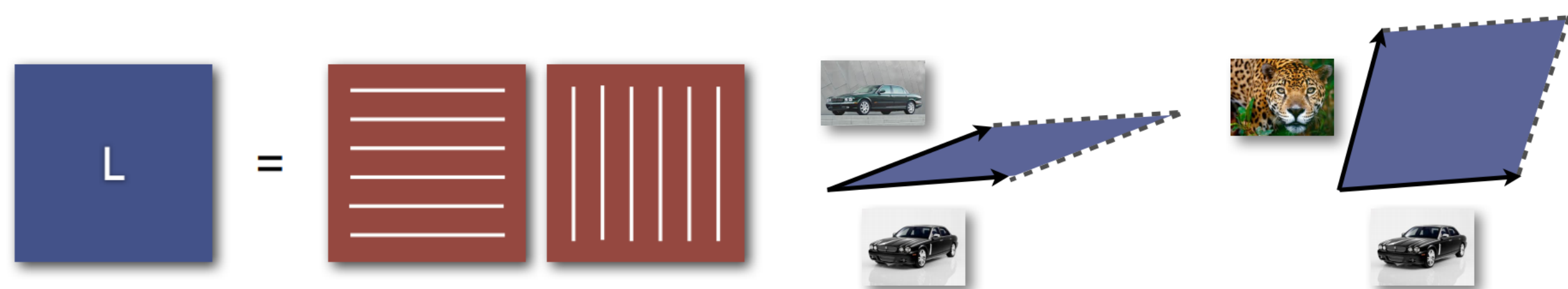
- $P(Y = \{i\}) = L_{ii}$
- Diagonal elements: always selected if close to 1

➤ Diversity of the subset

- $P(Y = \{i, j\}) \propto \begin{vmatrix} L_{ii} & L_{ij} \\ L_{ji} & L_{jj} \end{vmatrix} = L_{ii}L_{jj} - L_{ij}L_{ji}$
 $= P(i \in Y)P(j \in Y) - L_{ij}^2$
- Off-diagonal elements: Negative correlations
- Large Values of L_{ij} : $\{i, j\}$ less selected together

2) Geometric Intuition of DPP

- $L = B^T B$ where B is a $D \times N$ matrix (N : number of items)
- $P_L(Y) \propto \det(L_Y) = \text{Vol}^2(\{B_i\}_{i \in Y})$



DATASET & HYPERPARAMETERS

IWSLT14 German-English Dataset: translated TED talks

- 153,000 training | 7,000 development | 7,000 test pairs
- Byte Pair Encoding | Max Token 4096 | Max Seq 200

Key Hyper-parameters

- Epoch 300 | Embedding 512 | FC Embedding 1024
- # Attention Head 4 | # En/Decoder Layers 6

RESULTS & DISCUSSIONS

	Vanilla	Dis (A)	Dis (V)	Dis (O)	DPP (A)	DPP (V)	DPP (O)
BLUE	28.58				28.49	28.64	28.71
Time (h)	15		-		30.37	30.85	31.08

- In the future, additional experiments with larger dataset, like WMT14 will be needed.
- DPP can be applied to Transformer Networks with different structures, like Universal Transformer.
- Adaptive incorporation of quality and similarity score to build L-Matrix can be investigated as future work.