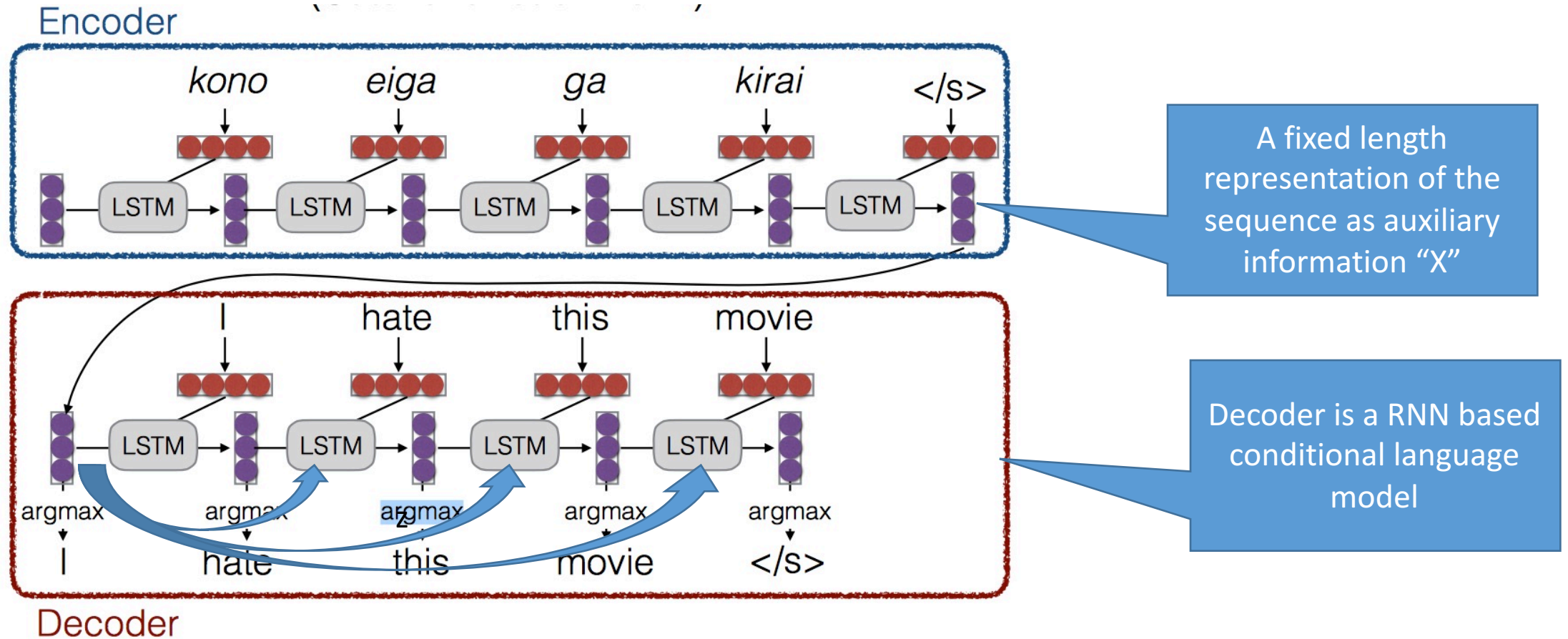# Attention in NLP

Presented by Yan Li

# Outline

- Basic seq2seq model & its drawback
- Basic attention and its application in NMT
- Variants of attention

# Seq2seq Models



Encoder

kono    eiga    ga    kirai    </s>

A fixed length representation of the sequence as auxiliary information "X"

Decoder

I    hate    this    movie

argmax    argmax    argmax    argmax    argmax

I    hate    this    movie    </s>

Decoder is a RNN based conditional language model

- Read whole sentience once and then translate it. Hard!
- "You can't cram the meaning of a whole sentence into a single vector!"— Ray Mooney
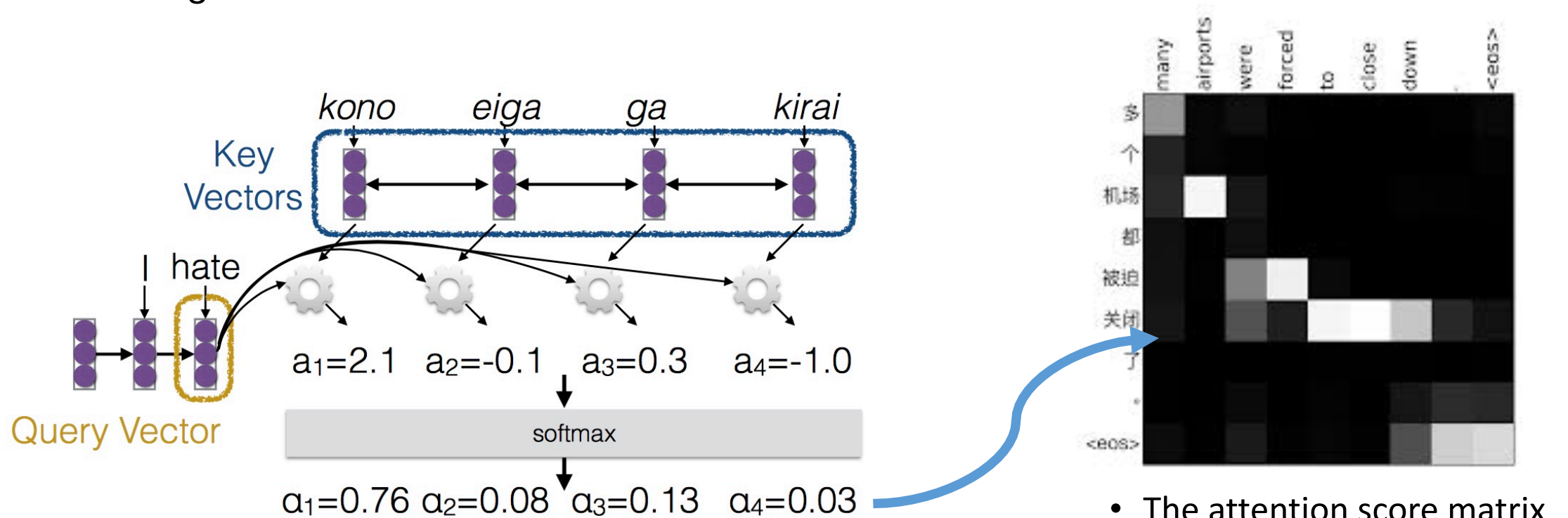
# Outline

- Basic seq2seq model & its drawback
- <span style="color:red">Basic attention and its application in NMT</span>
- Variants of attention

# Attention Basic Idea

- Encode each word in the sentence into a vector

- When decoding, perform a linear combination of these vectors, weighted by "attention weights"

- Use this combination in picking the next word

- Read the source sentence and when you translate you can go back to re-read the source sentence again and again. Relatively easy!
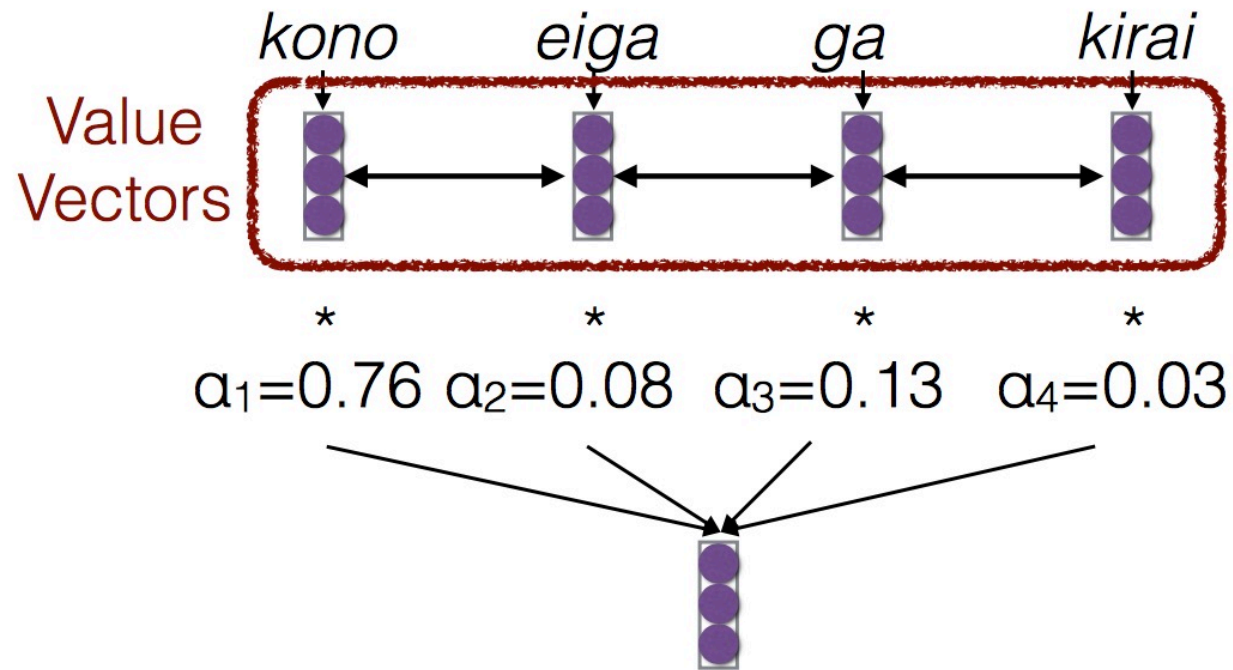
# Calculate Attention

1. Use "query" vector (decoder state) and "key" vectors (all encoder states)
2. For each query-key pair, calculate weight
3. Normalize to add to one using softmax
4. Combine together value vectors (usually encoder states, like key vectors) by taking the weighted sum

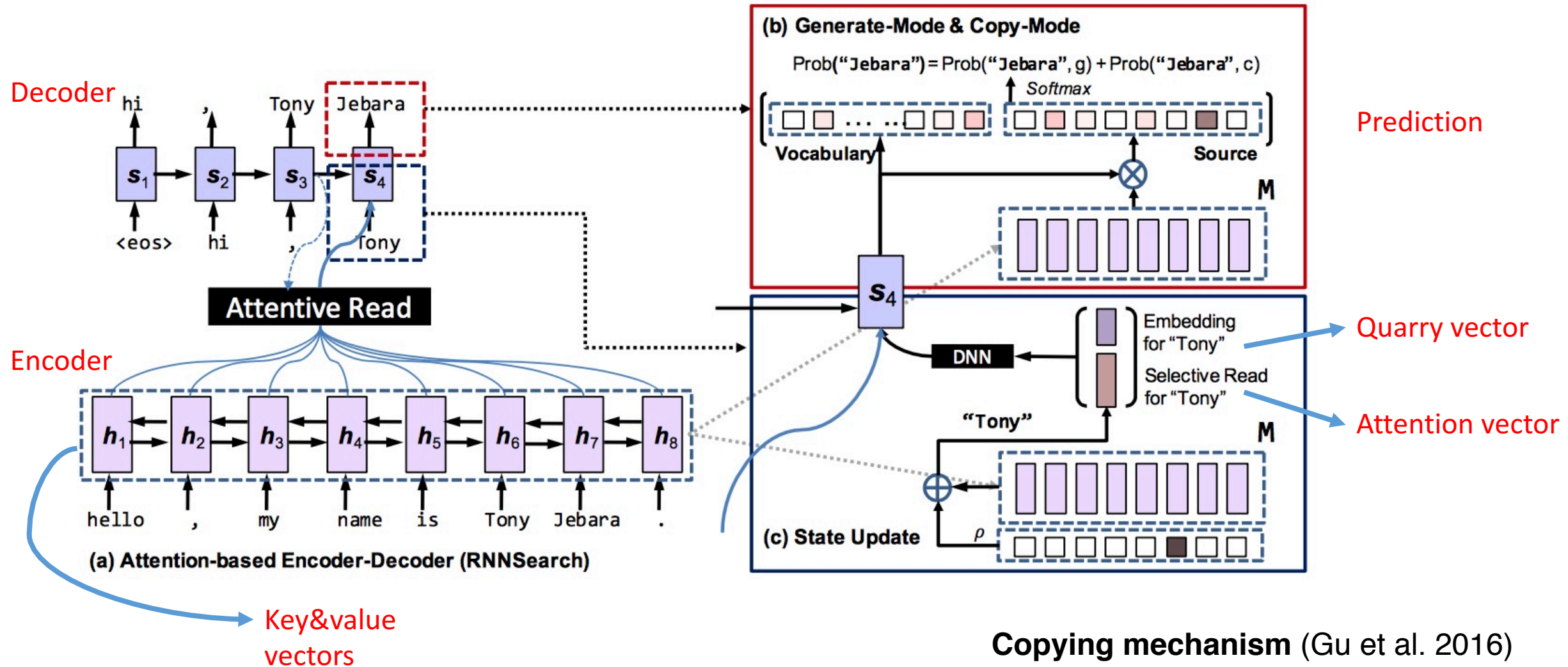

- The attention score matrix

# Calculating Attention (cont.)

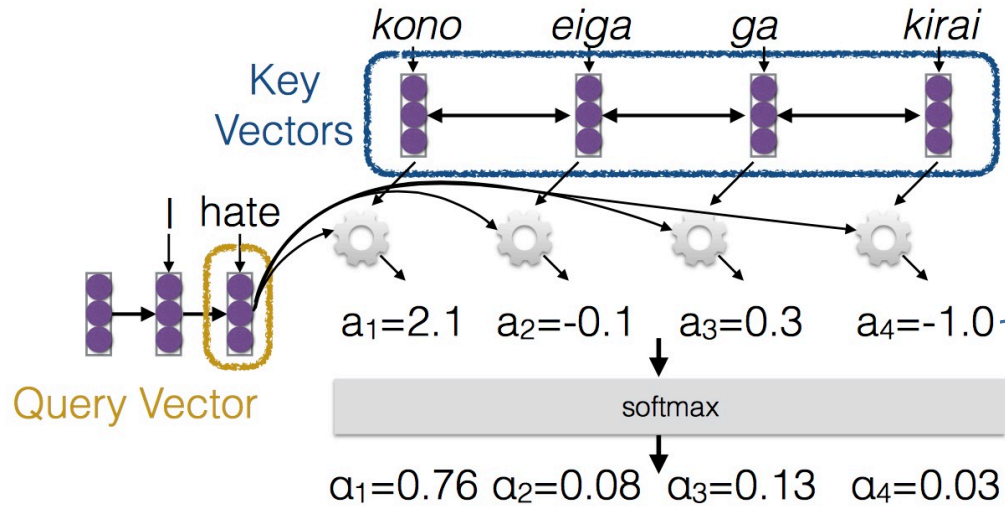- Combine together value vectors (usually encoder states, like key vectors) by taking the weighted sum

# Attend to the source language



(a) Attention-based Encoder-Decoder (RNNSearch)

(b) Generate-Mode & Copy-Mode

$$\text{Prob}(\text{"Jebara"}) = \text{Prob}(\text{"Jebara"}, g) + \text{Prob}(\text{"Jebara"}, c)$$

(c) State Update

**Copying mechanism** (Gu et al. 2016)

# Different Attention Score Functions



$q$ is the query and $k$ is the key

- **Multi-layer Perceptron:** $a(\boldsymbol{q}, \boldsymbol{k}) = \boldsymbol{w}_2^\mathsf{T} \tanh(W_1[\boldsymbol{q}; \boldsymbol{k}])$ (Flexible, often very good with large data)

- **Bilinear:** $a(\boldsymbol{q}, \boldsymbol{k}) = \boldsymbol{q}^\mathsf{T} W \boldsymbol{k}$

- **Dot Product:** $a(\boldsymbol{q}, \boldsymbol{k}) = \boldsymbol{q}^\mathsf{T} \boldsymbol{k}$ (No parameters, But requires sizes to be the same.)

- **Scaled Dot Product:** $a(\boldsymbol{q}, \boldsymbol{k}) = \dfrac{\boldsymbol{q}^\mathsf{T} \boldsymbol{k}}{\sqrt{|\boldsymbol{k}|}}$ (scale by size of the vector, this is because dot product increases as dimensions get larger)

# Outline

- Basic seq2seq model & its drawback
- Basic attention and its application in NMT
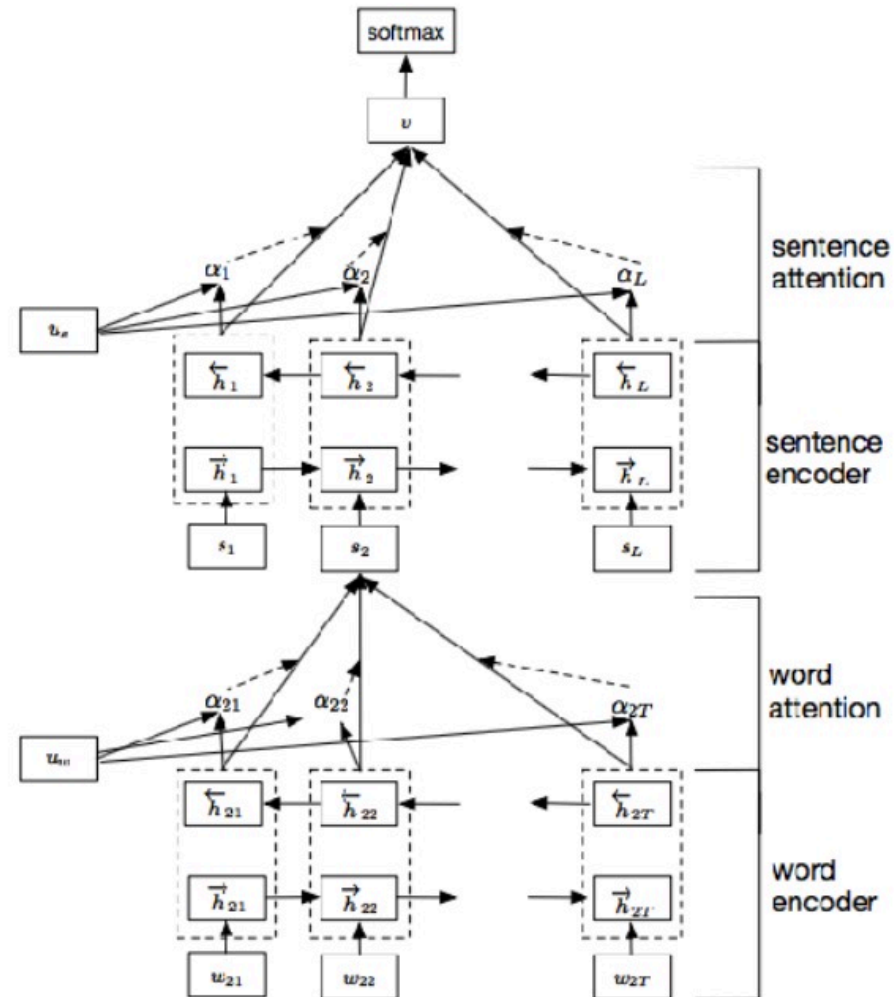- Variants of attention

# Attend to the previously generated things

- Key point: Read the previously generated words again and again while generate the new words.

- Reason: For example, in a long paragraph the key words might need many times.



$$p(\text{Yellen}) = g \, p_{\text{vocab}}(\text{Yellen}) + (1 - g) \, p_{\text{ptr}}(\text{Yellen})$$
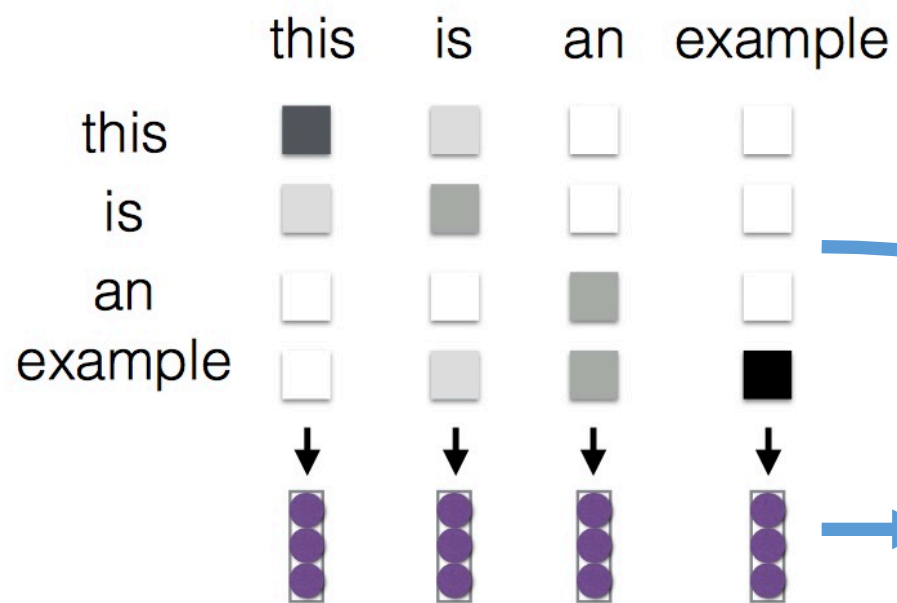
(Merity et al. 2016)

# Hierarchical Structures

- Encode with attention over each sentence, then attention over each sentence in the document. (Yang et al. 2016)

# Self Attention

- Each element in the sentence attends to other elements → context sensitive encodings! Incorporate the surrounding words information into the word embedding



$x_i$ is the original embedded matrix of i-th word

$$e_{ij} = \frac{(x_i W^Q)(x_j W^K)^T}{\sqrt{d_z}}$$

$$\alpha_{ij} = \frac{\exp e_{ij}}{\sum_{k=1}^{n} \exp e_{ik}}$$

$$z_i = \sum_{j=1}^{n} \alpha_{ij}(x_j W^V)$$

- Comparing with Bi-RNN:
  - Much faster than Bi-RNN.
  - More direct than Bi-RNN

Each output element, $z_i$, is computed as weighted sum of a linearly transformed input elements

# Attention Is All You Need (Vaswani et al. 2017)

(The probability of next word)

Output Probabilities

Softmax

Linear

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Nx

Add & Norm

Feed Forward

Nx

Add & Norm

Multi-Head Attention

Add & Norm

Masked Multi-Head Attention

Positional Encoding

Input Embedding

Inputs

Positional Encoding

Output Embedding

Outputs (shifted right)

**The above slides mentioned attention**

**Self-Attention over input sentence (encode source)**

**Make sure that even if no RNN, it can still distinguish positions**

**Encode the previous generated words**

(The target we produce so far)

Every step is one training sample, in RNN the whole sentence is one sample, you have to backprob all sentence. Here each token is a sample, one don't need to backprob whole sentence.

Multi-Head Attention

Linear

Concat

Scaled Dot-Product Attention
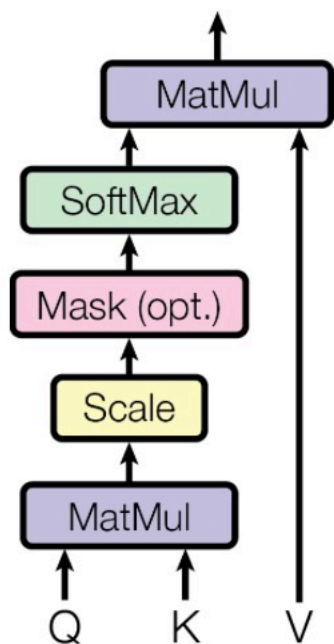
h

Linear  Linear  Linear

V  K  Q

**From Encoding**   **From Decoding**

This structure is used several times. In the arrowhead one the input from different part

# Important components of Attention Is All You Need
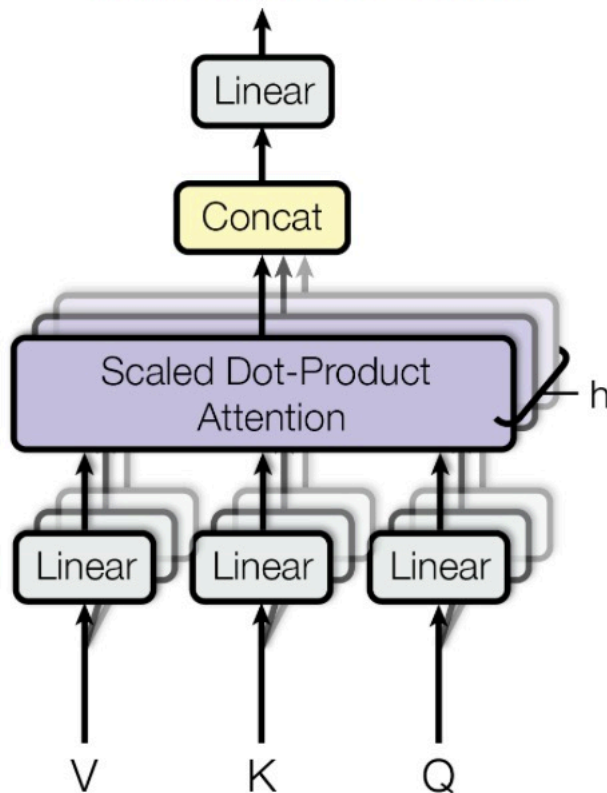
## Scaled Dot-Product Attention



$$\mathrm{Attention}(Q, K, V) = \mathrm{softmax}(\frac{QK^T}{\sqrt{d_k}})V$$

Q: query vector
K: Key vector
V: value vector

## Multi-Head Attention



constitute

**Idea:** multiple attention "heads" focus on different parts of the sentence

- Positional Encoding

use sine and cosine functions of different frequencies:

$$PE_{(pos,2i)} = sin(pos/10000^{2i/d_{\mathrm{model}}})$$

$$PE_{(pos,2i+1)} = cos(pos/10000^{2i/d_{\mathrm{model}}})$$

where $pos$ is the position and $i$ is the dimension.

感觉和**傅立叶变换**有相似的地方，把绝对位置表示为相对位置。