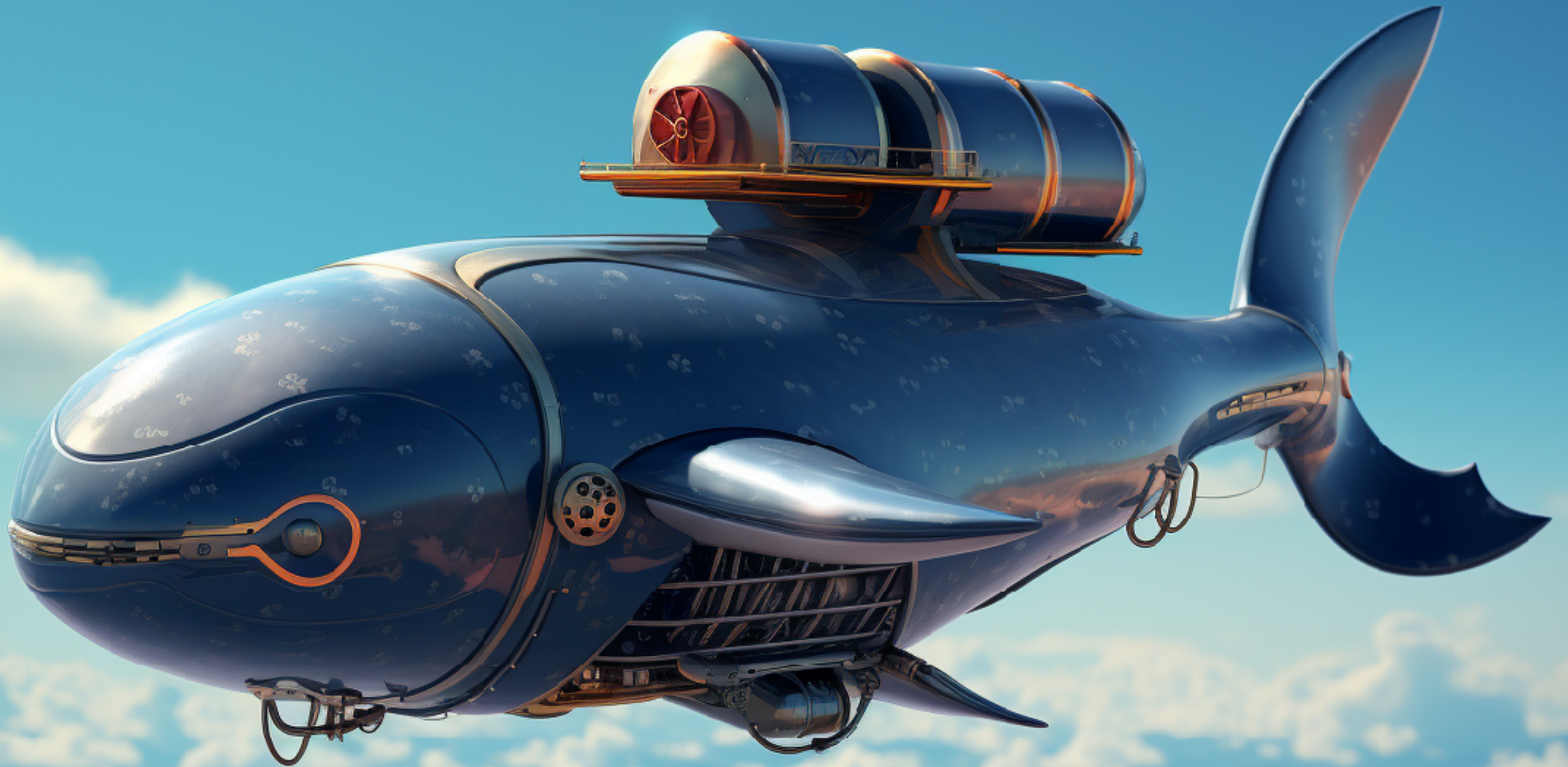


Transformerを詳しく学ぼう！



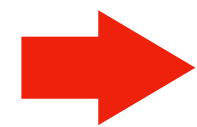
Section2

Section2の概要



講座の内容

Section1. Transformerの概要



Section2. Attentionの仕組み

Section3. Transformerにおける埋め込み

Section4. Transformerを組み立てる

今回の内容

1. Section2の概要
2. Attentionの概要
3. Scaled Dot-Product Attention
4. Multi-HeadAttention
5. 演習

教材の紹介

- **Pythonの基礎:**

python_basic

- **Section2の教材**

01_scaled_dot_product_attention.ipynb

02_multi_head_attention.ipynb

03_exercise.ipynb

https://github.com/yukinaga/learning_transformer/

Section1演習の解答例

- 03_exercise.ipynb

Attentionの概要



ChatGPT に聞いてみる

「Transformerで使われるAttentionって何ですか？」

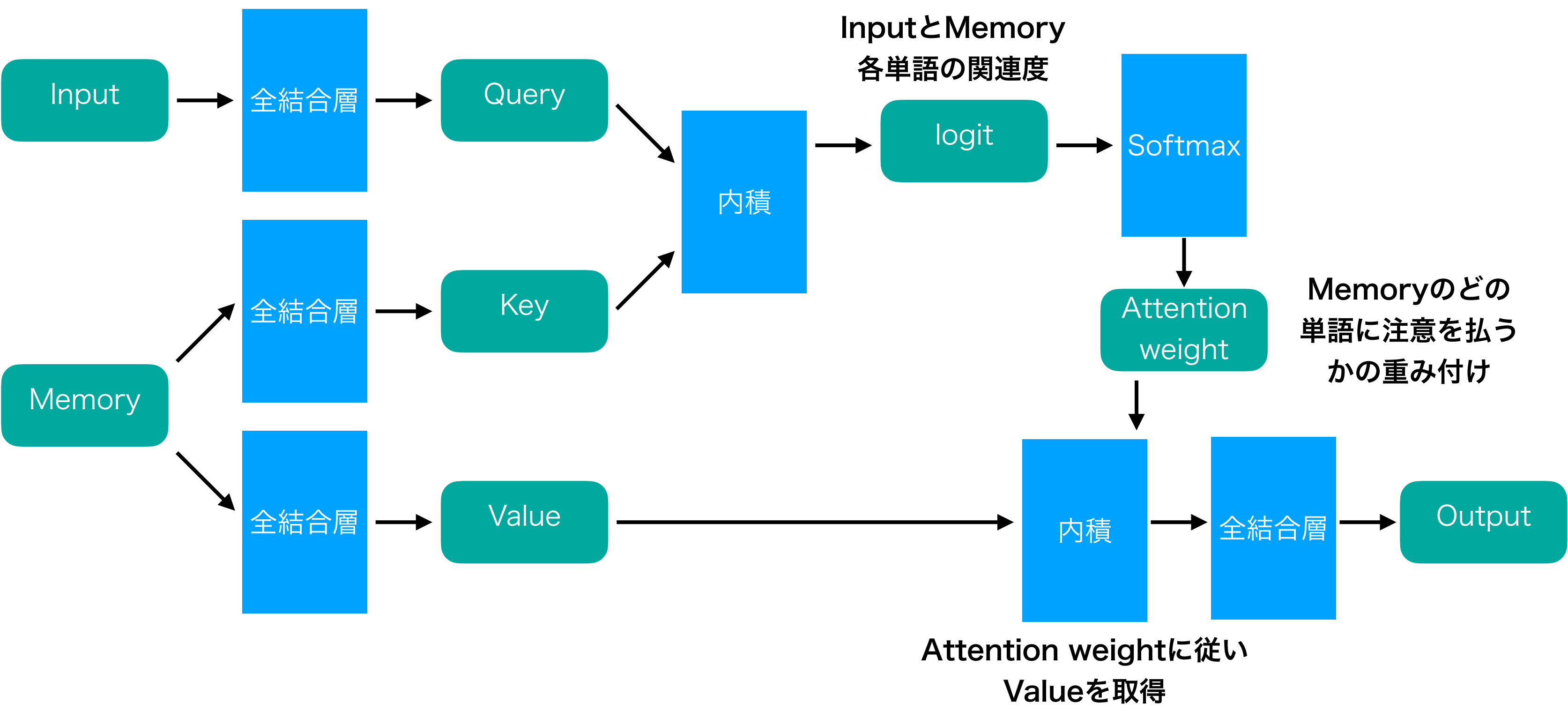
→ ?

<https://chat.openai.com/>

Attentionとは？

- **Attentionとは？**
 - 文章中のどの単語に注目すればいいかを表すスコア
 - Query、Key、Valueの3つのベクトルで計算される
- **Query**
 - Inputのうち「検索をかけたいもの」
- **Key**
 - 検索対象とQueryの近さを測る
- **Value**
 - Keyに基づき、適切なValueを出力する

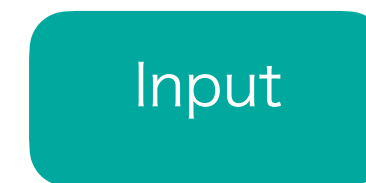
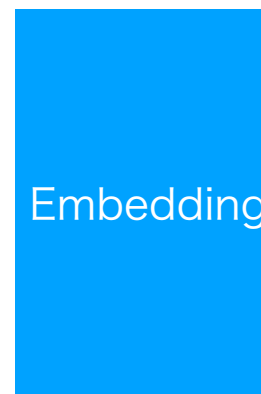
Attentionとは？



InputとMemory

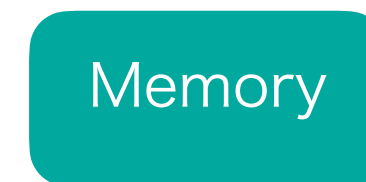
得意 な スポーツ は ？

(各単語はidで表される)

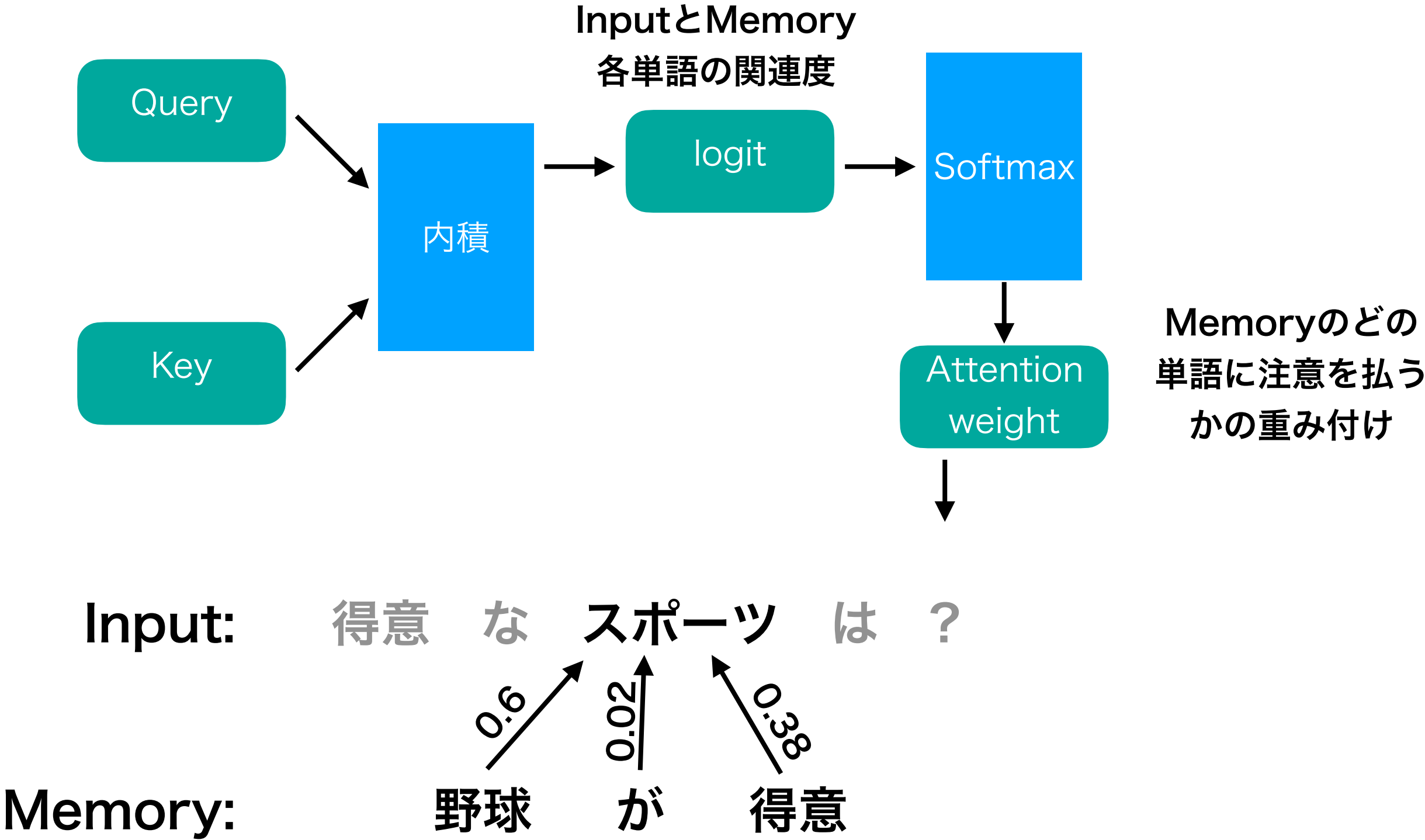


野球 が 得意

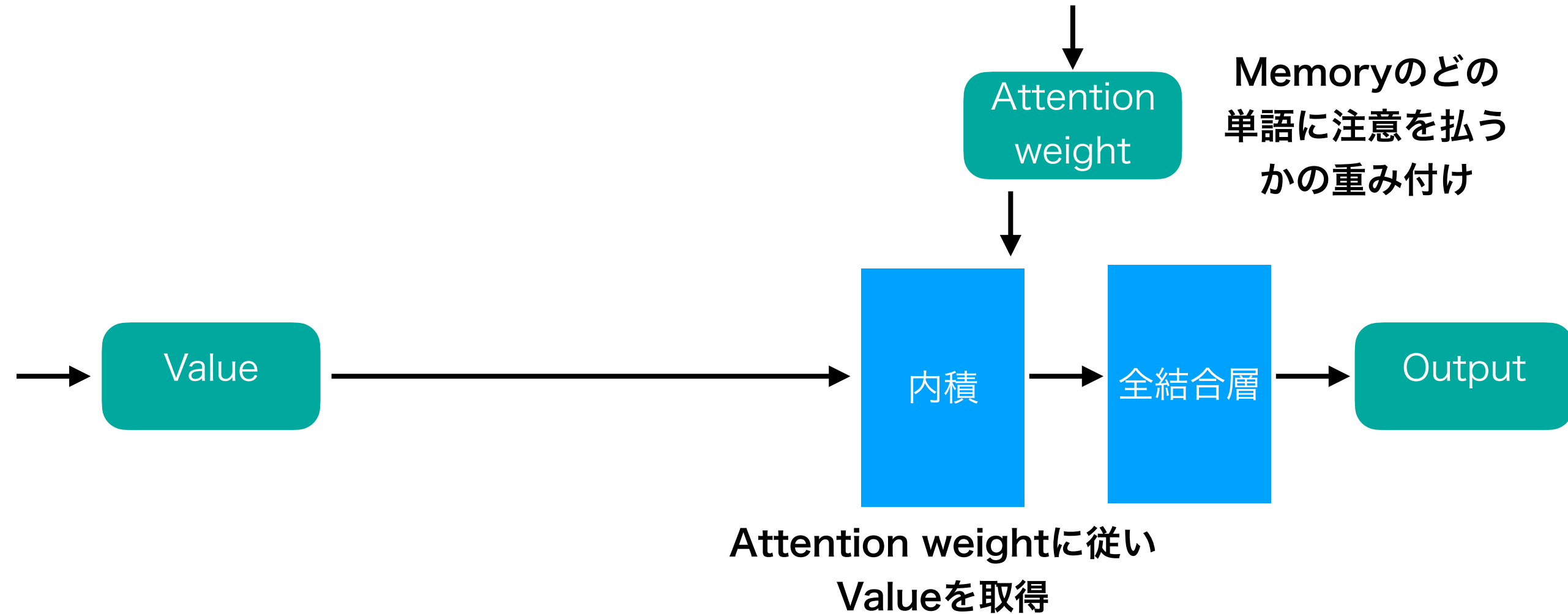
(各単語はidで表される)



Attention weightの計算



Valueと内積



Input: 得意 な スポーツ は ?

Memory:

野球 が 得意

0.6 ↗
0.02 ↑
0.38 ↘

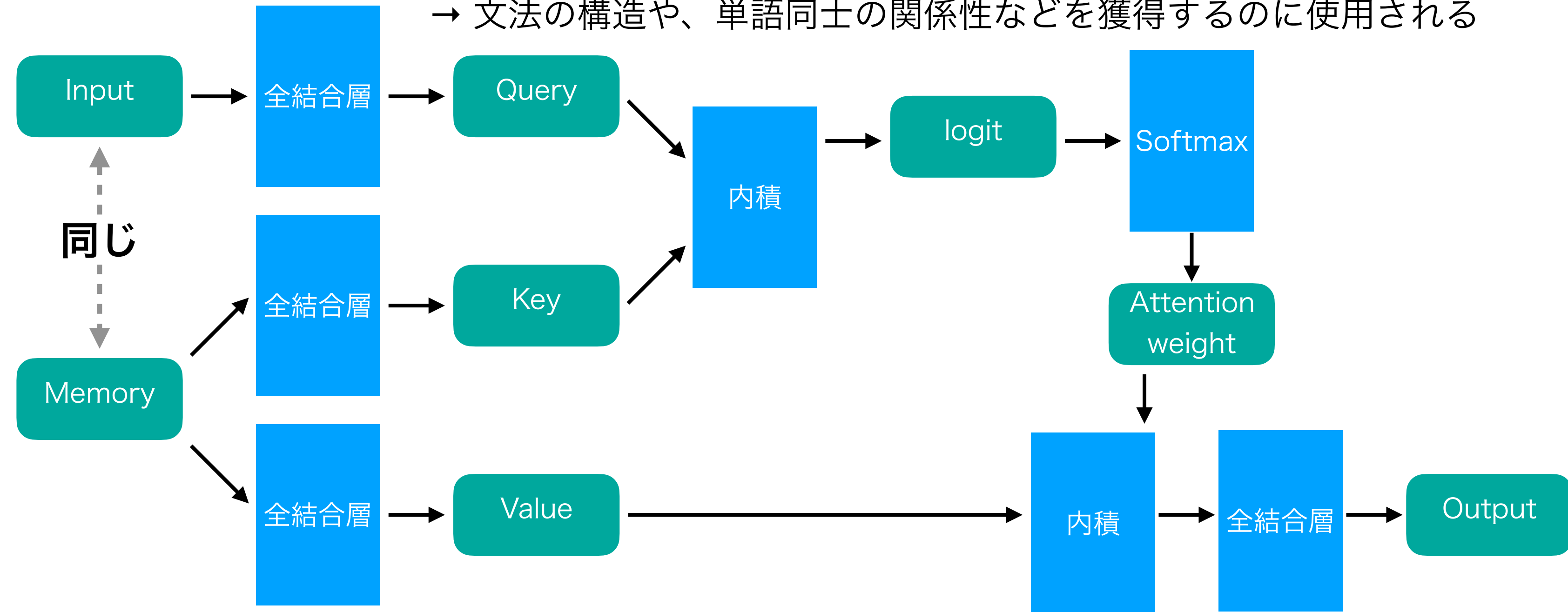
$$\begin{aligned}\text{内積} &= \text{Value}(\text{野球}) \times 0.6 \\ &\quad + \text{Value}(\text{が}) \times 0.02 \\ &\quad + \text{Value}(\text{得意}) \times 0.38\end{aligned}$$

Self-Attention

- **Self-Attention**

→ InputとMemoryが同一のAttention

→ 文法の構造や、単語同士の関係性などを獲得するのに使用される

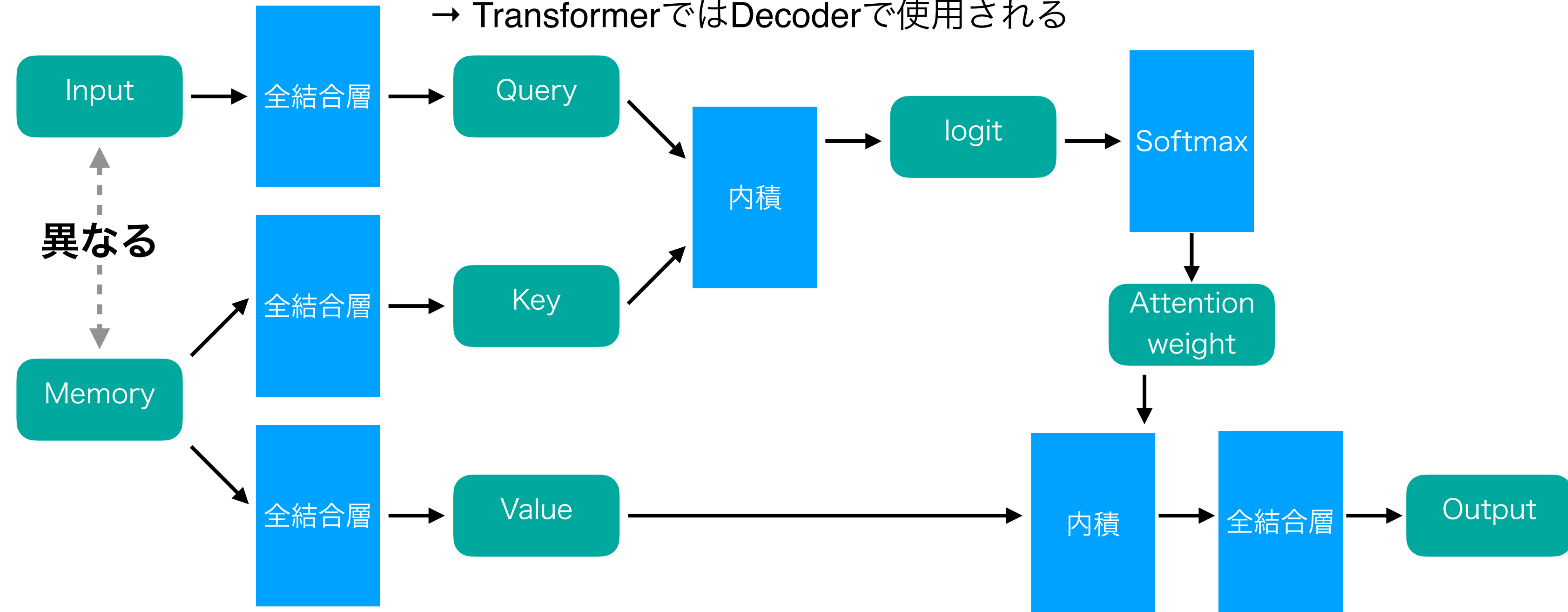


SourceTarget-Attention

- **Self-Attention**

→ InputとMemoryが異なるAttention

→ TransformerではDecoderで使用される



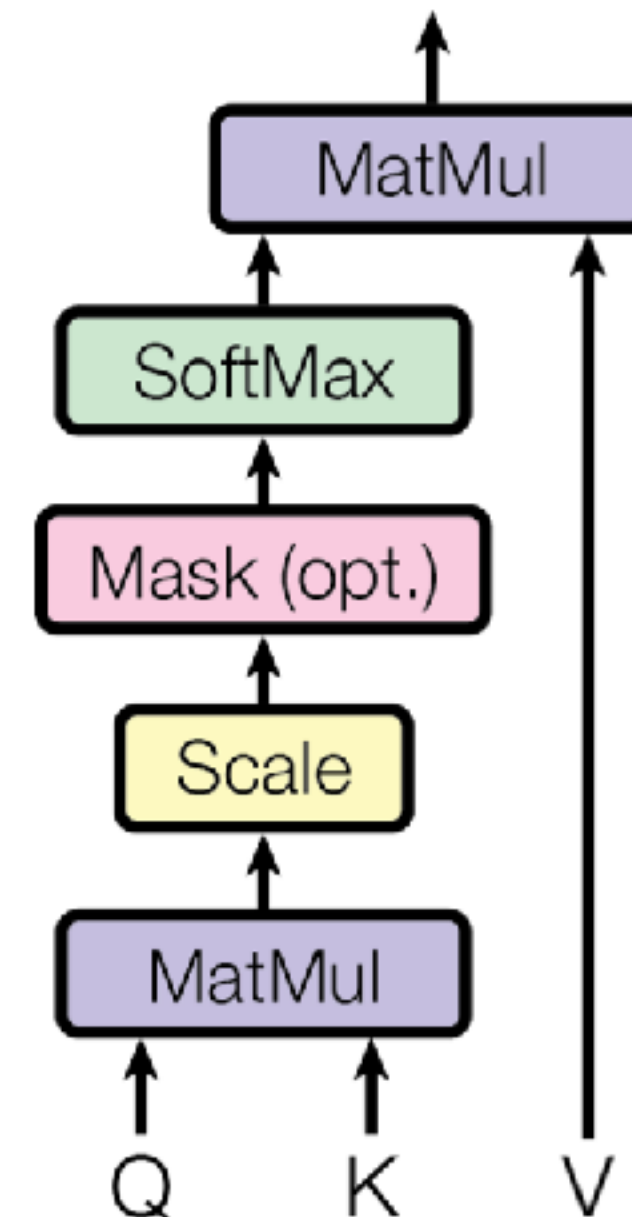
Scaled Dot-Product Attention

- **Scaled Dot-Product Attention**

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

- **Mask**

→ 入力した単語が「先読み」されるのを防ぐために、特定の key に対して
Attention weight を0にする

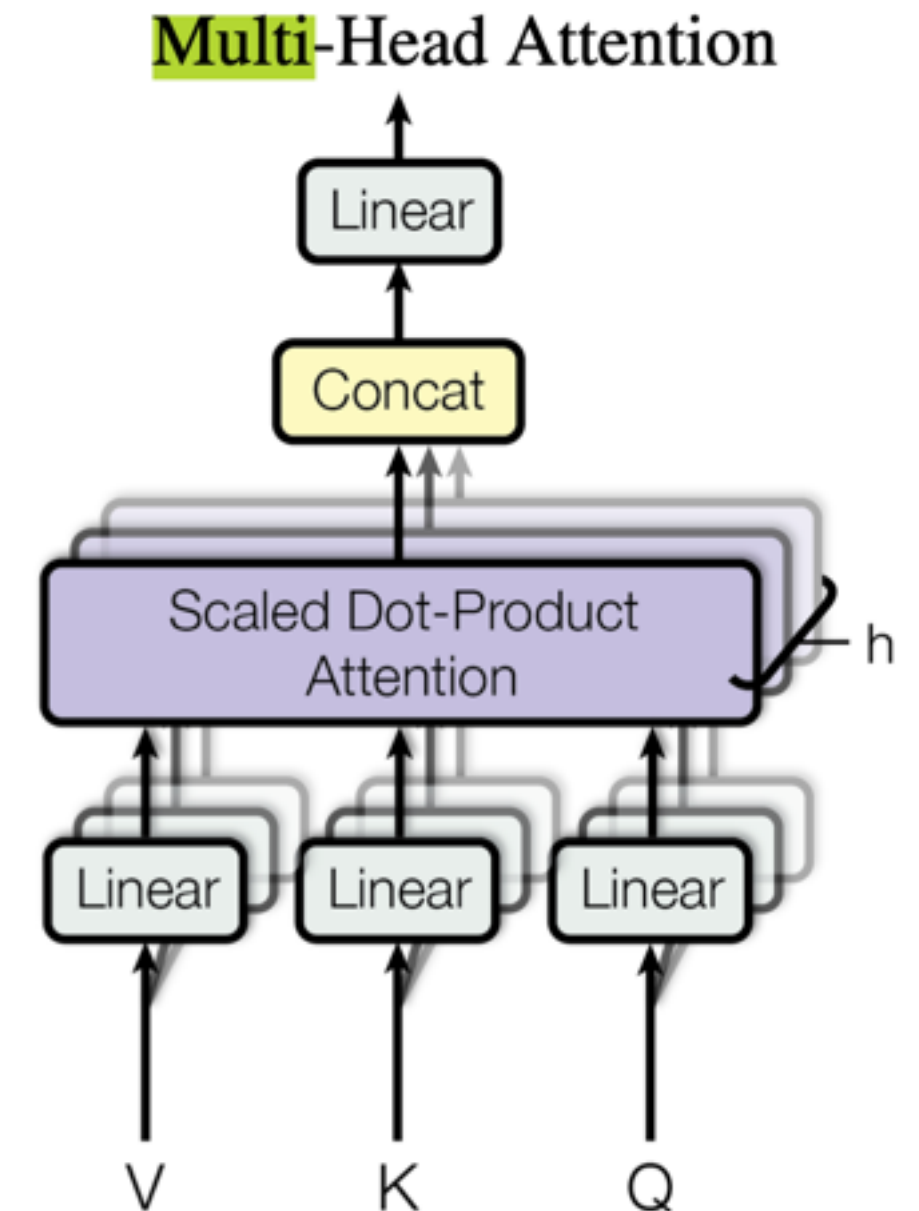


Multi-Head Attention

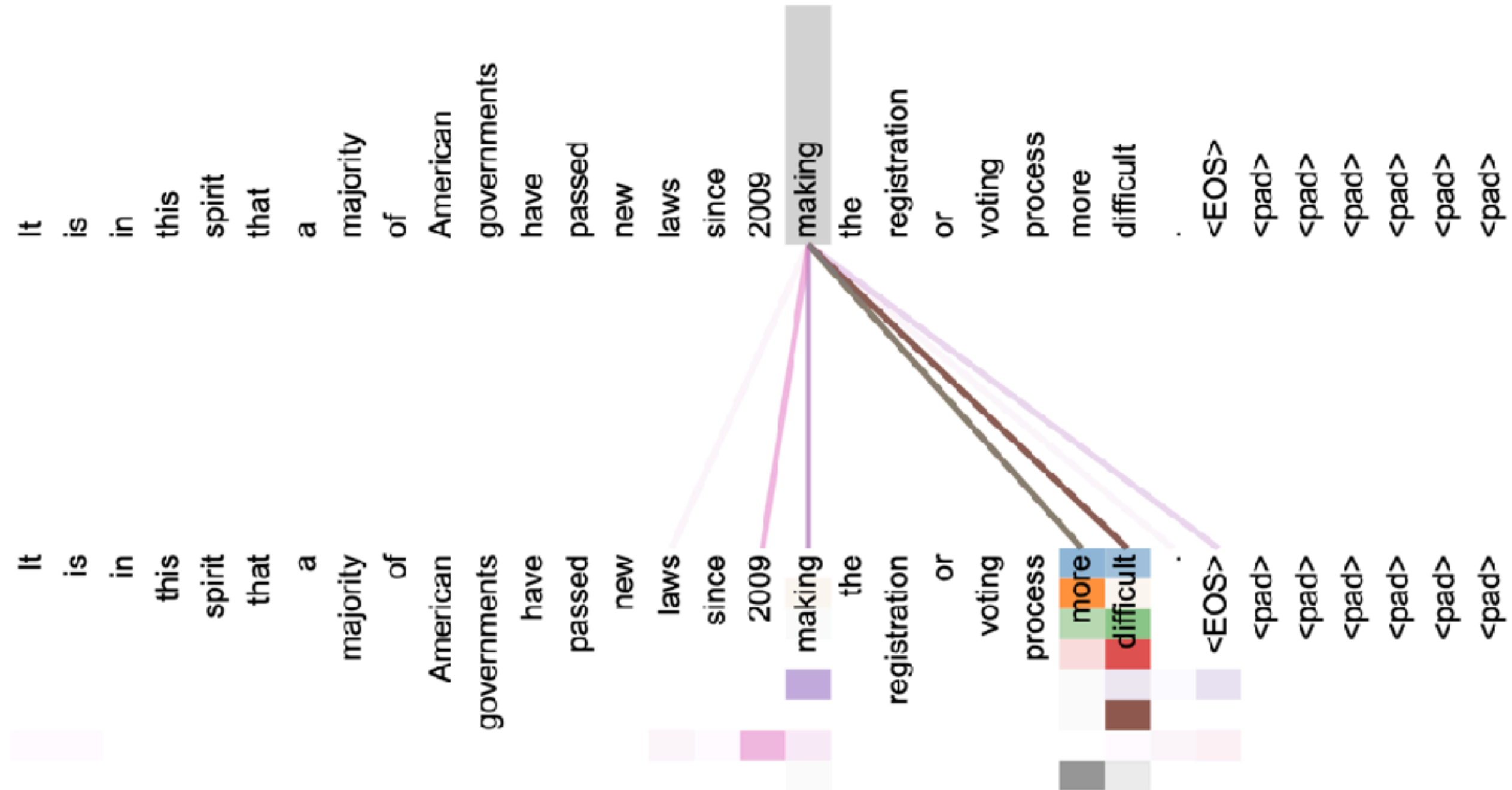
- **Multi-Head Attention**
 - Attentionを並行に並べる
 - それぞれのAttentionはHeadと呼ばれる
 - 「Attention Is All You Need」ではMulti-Head化による性能の向上が述べられている

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h) W^O$$

where $\text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$



Attentionの可視化



異なる色は異なるAttentionのHeadを表す

Attention Is All You Need, Ashish, V. et al. (2017) より引用



休憩...

Scaled Dot-Product Attention



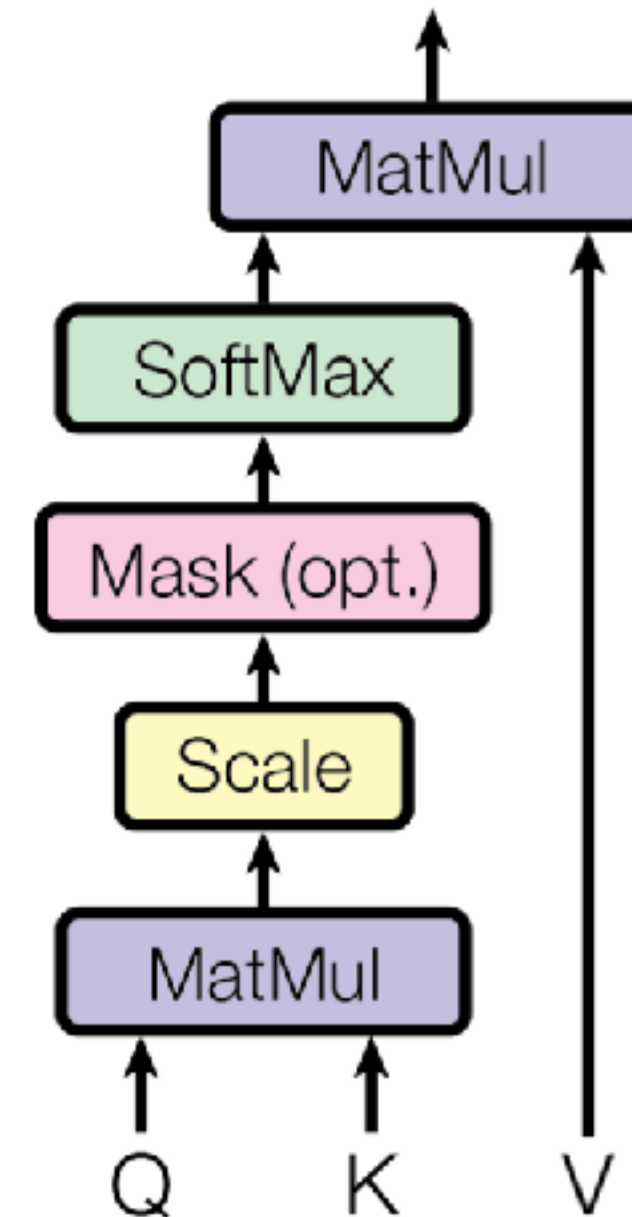
Scaled Dot-Product Attention

- **Scaled Dot-Product Attention**

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

- **Mask**

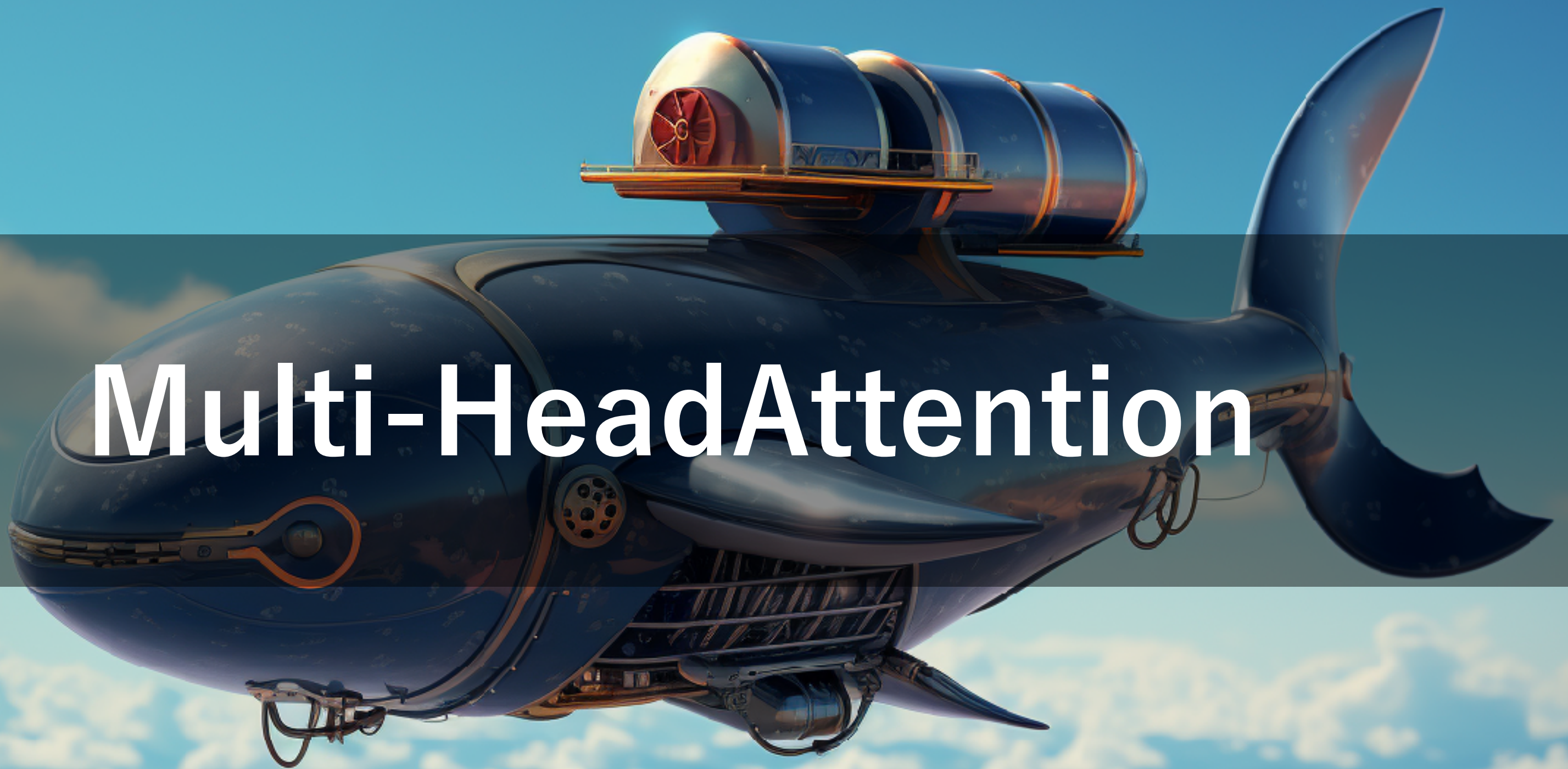
→ 入力した単語が「先読み」されるのを防ぐために、特定の key に対して
Attention weight を0にする



Scaled Dot-Product Attention

- 01_scaled_dot_product_attention.ipynb

Multi-HeadAttention

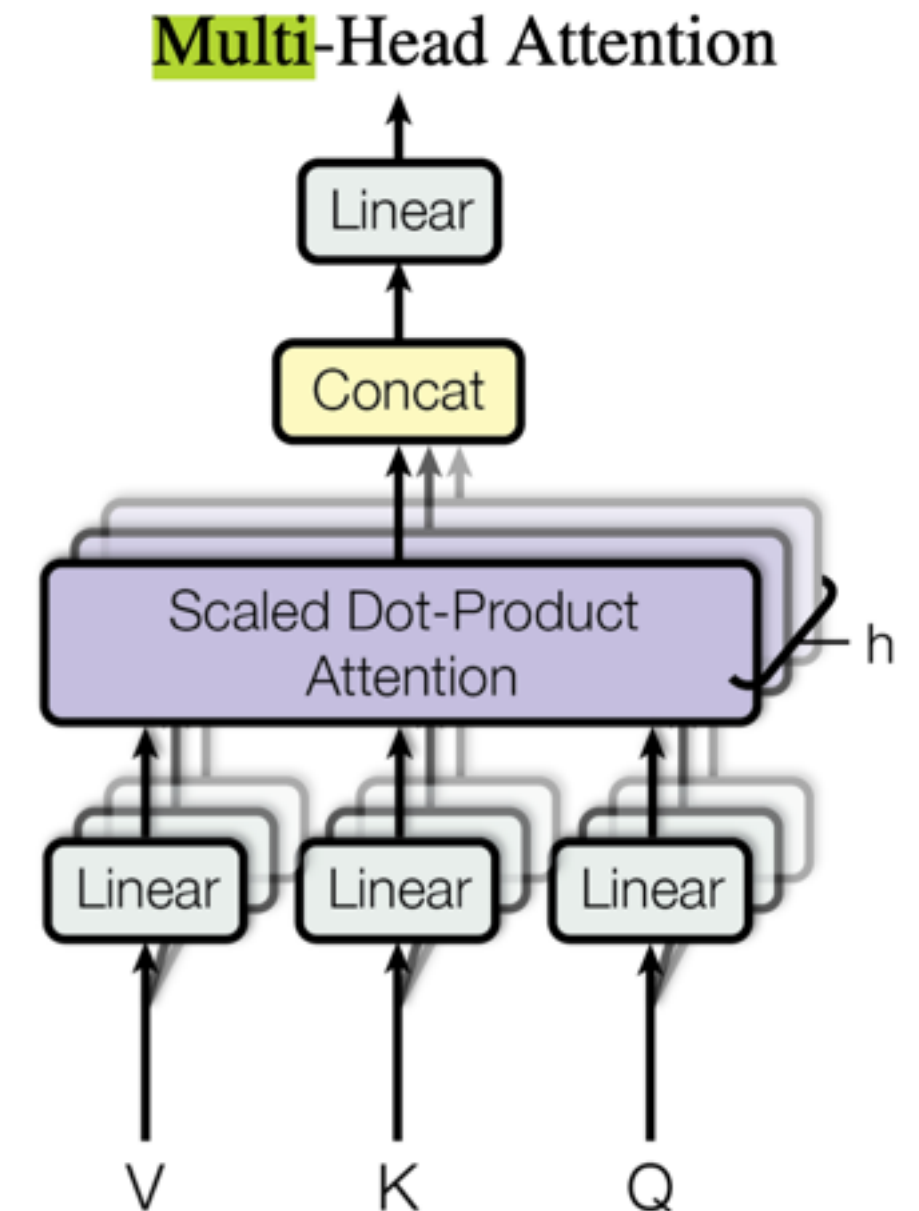


Multi-Head Attention

- **Multi-Head Attention**
 - Attentionを並行に並べる
 - それぞれのAttentionはHeadと呼ばれる
 - 「Attention Is All You Need」ではMulti-Head化による性能の向上が述べられている

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h) W^O$$

where $\text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$



Multi-HeadAttention



- 02_multi_head_attention.ipynb

演習



演習

- 03_exercise.ipynb

次回の内容

Section1. Transformerの概要

Section2. Attentionの仕組み

 **Section3. Transformerにおける埋め込み**

Section4. Transformerを組み立てる