

# 注意機構とトランスフォーマー

本谷 秀堅

名古屋工業大学

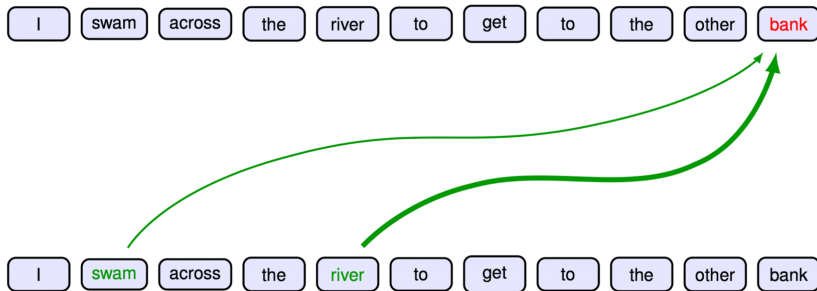
# Transformer (トランスフォーマー)

注意機構 (Attention) を核とするニューラルネットワーク

# 言語データ（再掲）

## 言語データ

- シンボルの系列
- 離れたシンボル間に強い関係
- シンボル間の非対称な関係（例：対義語・類義語・包含関係）



# Word2Vec による単語ベクトルの埋め込み

- 埋め込み先は文脈に依存せずに同じ
  - bank を表すベクトルは、堤のときも銀行のときも同じ

$$\mathbf{h} = \mathbf{W}\mathbf{x}$$

## 分布仮説

単語の意味はその単語の周辺に現れる単語、すなわち文脈 (context) によって決まる

# 注意機構の目的

各単語の変換先を、系列中の他の単語に依存して変化させたい

- I swam across the river to get to the other **bank**.
  - “bank” は river の近くへと変換
- I walked across the road to get cash from the **bank**.
  - “bank” は money の近くへと変換

注意機構により文脈に依存して各単語の埋め込み先を変化させる

- 特徴ベクトルのことを「トークン (token)」と呼ぶ。

文脈中のすべての単語を参照して、各単語の埋め込み先を決める

- 入力： $N$  単語。それぞれ  $D$  次元のベクトル（トークン）で表現

$$\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N \quad (\mathbf{x}_n \in \mathbb{R}^D)$$

- 出力：各単語の変換後の  $D$  次元トークン

$$\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_N \quad (\mathbf{y}_n \in \mathbb{R}^D)$$

文脈中の各単語の荷重平均を埋め込み先とする

$$\mathbf{y}_n = \sum_{m=1}^N a_{nm} \mathbf{x}_m$$

# 注意プーリング (Attention Pooling)

注意プーリング (pooling: 集約するの意)

$a_{nm}$ : 注意の強さ ( $a_{nm} > 0$ ,  $\sum_{m=1}^N a_{nm} = 1$ )

$$\mathbf{y}_n = \sum_{m=1}^N a_{nm} \mathbf{x}_m$$

注意の強さを文脈に応じて決めたい

- I swam across the river to get to the other **bank**.
  - “bank” は river と似た方向へと変換
- I walked across the road to get cash from the **bank**.
  - “bank” は money と似た方向へと変換

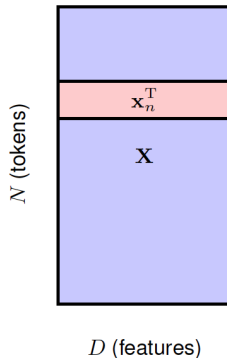


- $N$  単語
- 各単語はベクトルで表現  $\mathbf{x}_n \in \mathbb{R}^D$
- 文脈は行列で表現:  $\mathbf{X} \in \mathbb{R}^{N \times D}$

$$\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n]^\top$$

- 出力も行列で表現:  $\mathbf{Y} \in \mathbb{R}^{N \times D}$

$$\mathbf{Y} = [\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_n]^\top$$



# 注意の大きさを文脈で決める (案1)

もともと似たベクトルどうしを強い重みで

$$a_{nm} = \frac{\exp(\mathbf{x}_n^\top \mathbf{x}_m)}{\sum_k \exp(\mathbf{x}_n^\top \mathbf{x}_k)}$$

↓

$$\mathbf{A} = \text{SM}(\mathbf{X}\mathbf{X}^\top) \in \mathbb{R}^{N \times N},$$

$$\mathbf{Y} = \mathbf{A}\mathbf{X} = \text{SM}(\mathbf{X}\mathbf{X}^\top) \in \mathbb{R}^{N \times N}.$$

- トークン間の類似度を内積で評価する場合の式
- 行列  $\mathbf{A}$  の  $(n, m)$  要素が  $a_{nm}$
- Softmax は各行で計算
- $\mathbf{X}\mathbf{X}^\top$  は対称行列
- $a_{nm}$  に「学習」の余地無し。柔軟性に欠ける。

## 注意の大きさを文脈で決める (案2)

注意の大きさの決め方を学習する

$$\tilde{\mathbf{X}} = \mathbf{X}\mathbf{U}$$

上記  $\mathbf{U} \in \mathbb{R}^{D \times D}$  を学習する

$$\mathbf{A} = \text{SM}(\mathbf{X}\mathbf{U}\mathbf{U}^\top \mathbf{X}^\top),$$

$$\mathbf{Y} = \text{SM}(\mathbf{X}\mathbf{U}\mathbf{U}^\top \mathbf{X}^\top) \mathbf{X}\mathbf{U}.$$

- 案1より柔軟
- $\mathbf{X}\mathbf{U}\mathbf{U}^\top \mathbf{X}^\top$  は対称行列
- Softmax関数を適用すると対称ではなくなるが大小関係は維持される
- 包含関係、対義語、類義語、指示代名詞など複雑な関係の表現には自由度が不足

# 注意の大きさを文脈で決める（注意機構）

注意の大きさの決め方：クエリ・キー・バリューの考え方を採用する

$$a_{nm} = \frac{\exp(\mathbf{q}_n^\top \mathbf{k}_m)}{\sum_k \exp(\mathbf{q}_n^\top \mathbf{k}_k)}$$

↑

$$\mathbf{Q} = \mathbf{X}\mathbf{W}^{(q)},$$

$$\mathbf{K} = \mathbf{X}\mathbf{W}^{(k)},$$

$$\mathbf{V} = \mathbf{X}\mathbf{W}^{(v)}.$$

行列  $\mathbf{W}^{(q)}$ ,  $\mathbf{W}^{(k)}$ ,  $\mathbf{W}^{(v)}$  を学習により決める。

$\mathbf{W}^{(q)}, \mathbf{W}^{(k)} \in \mathbb{R}^{D \times D_k}$ ,  $\mathbf{W}^{(v)} \in \mathbb{R}^{D \times D_v}$

$$\mathbf{Y} = \text{SM}(\mathbf{Q}\mathbf{K}^\top)\mathbf{V}$$

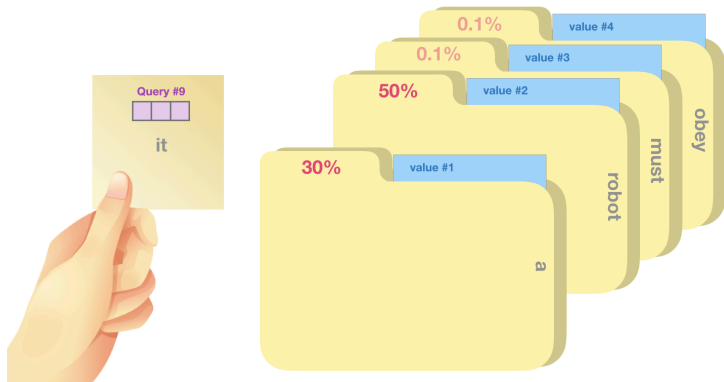
# (復習) クエリ・キー・バリュー

- キー (key) とバリュー (value: 値) の組のデータ集合
- クエリ (query) による検索

例：

- データ集合:  
 $\mathcal{D} = \{(\text{“夏目”}, \text{“漱石”}), (\text{“森”}, \text{“鷗外”}), (\text{“中島”}, \text{“敦”}), (\text{“松尾”}, \text{“芭蕉”})\}$
- 検索：query = “森”  
検索結果：value = “鷗外”

# (復習) クエリ・キー・バリュー



$$Y = \text{SM}(\mathbf{Q}\mathbf{K}^\top)\mathbf{V}$$

- ① 行列  $\mathbf{W}^{(q)}$ ,  $\mathbf{W}^{(k)}$  を学習する

- key-vector の次元を  $D_k$  で表すと、いずれも  $\mathbb{R}^{D \times D_k}$

- ②  $\mathbf{Q} = \mathbf{X}\mathbf{W}^{(q)}$ ,  $\mathbf{K} = \mathbf{X}\mathbf{W}^{(k)}$

- $n$  番目の単語のクエリ ( $\mathbf{q}_n$ ) とキー ( $\mathbf{k}_n$ ) は次のとおり

$$\mathbf{q}_n^\top = \mathbf{x}_n^\top \mathbf{W}^{(q)}, \quad \mathbf{k}_n^\top = \mathbf{x}_n^\top \mathbf{W}^{(k)}$$

- ③  $n$  番目の単語を変換するために  $m = 1, 2, \dots, N$  番目の単語を参照

$$\mathbf{q}_n^\top \mathbf{k}_m = (\mathbf{Q}\mathbf{K}^\top)_{nm}$$

- ④  $\mathbf{Q}\mathbf{K}^\top$  の各行で Softmax を計算

$$a_{nm} = \frac{\exp(\mathbf{q}_n^\top \mathbf{k}_m)}{\sum_k \exp(\mathbf{q}_n^\top \mathbf{k}_k)}$$

$$\mathbf{Y} = \text{SM}(\mathbf{Q}\mathbf{K}^\top)\mathbf{V}$$

- 行列  $\mathbf{W}^{(v)} \in \mathbb{R}^{D \times D}$  を学習する
- バリュウ (value) の計算:  $\mathbf{V} = \mathbf{X}\mathbf{W}^{(v)}$ 
  - $n$  番目の単語のバリュウ,  $\mathbf{v}_n$

$$\mathbf{v}_n^\top = \mathbf{x}_n^\top \mathbf{W}^{(v)}$$

- トークンの更新
  - $n$  番目の単語更新用の注意:  $\mathbf{a}_n^\top = [a_{n1}, a_{n2}, \dots, a_{nN}]$
  - $n$  番目の単語のトークンの更新

$$\mathbf{y}_n^\top = \mathbf{a}_n^\top \mathbf{V}$$



# 注意のスケールの調整

$\exp$  中の絶対値が大きくなることを防ぐ。(絶対値が大きいと勾配がゼロに近づき学習に失敗する。)

$$a_{nm} = \frac{\exp(\mathbf{q}_n^\top \mathbf{k}_m)}{\sum_k \exp(\mathbf{q}_n^\top \mathbf{k}_k)}$$

↓

$$a_{nm} = \frac{\exp\left(\frac{\mathbf{q}_n^\top \mathbf{k}_m}{\sqrt{D_k}}\right)}{\sum_k \exp\left(\frac{\mathbf{q}_n^\top \mathbf{k}_k}{\sqrt{D_k}}\right)}$$

$D_k$  は  $\mathbf{q}$ ,  $\mathbf{k}$  の次元

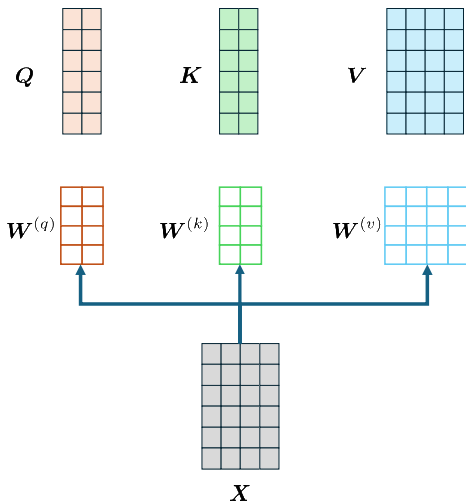
$$\mathbf{Y} = \text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{SM} \left( \frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{D_k}} \right) \mathbf{V}$$

# Attention Head

文脈に依存して埋め込み先を変える機構：注意機構

Attention Head と呼ぶ

- 入力文全体を見渡して「どこに注目するか」を判断するモジュール

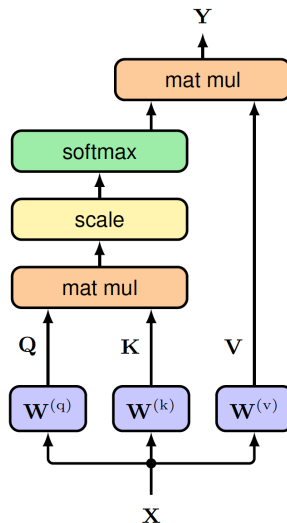


# Attention Head

文脈に依存して埋め込み先を変える機構：注意機構

Attention Head と呼ぶ

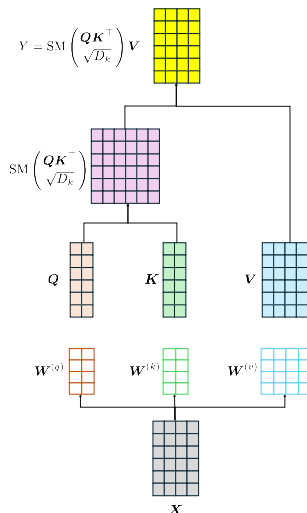
- 入力文全体を見渡して「どこに注目するか」を判断するモジュール



文脈に依存して埋め込み先を変える機構：注意機構

Attention Head と呼ぶ

- 入力文全体を見渡して「どこに注目するか」を判断するモジュール

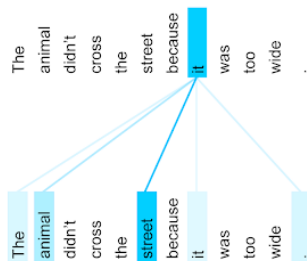
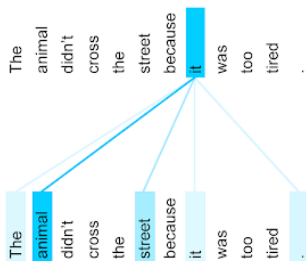


# Attention Head で実現できたこと

- 文脈に依存して単語の埋め込み先を変える
- 注意の強弱を文脈に依存して決定

例

The animal did not cross the street because it was too tired.  
The animal did not cross the street because it was too wide.



# 注意機構の性質

重要： $\mathbf{W}^{(q)}$ ,  $\mathbf{W}^{(k)}$ ,  $\mathbf{W}^{(v)}$  が注意の規準を定める（学習後は規準固定）

- 出力  $\mathbf{y}_n$  は入力  $\mathbf{v}_n$  の線形和

$$\mathbf{y}_n = \sum_{m=1}^N a_{nm} \mathbf{v}_m$$

- $\mathbf{v}_n$  は  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$  の線形変換

$$\mathbf{V} = \mathbf{XW}^{(o)}$$

- $\mathbf{y}_n$  は入力  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$  の線形変換
- 非線形性は注意の強度,  $a_{nm}$ , 経由で導入 ( $\mathbf{x}_n$  に依存・Softmax)

# 注意機構の性質

重要： $\mathbf{W}^{(q)}$ ,  $\mathbf{W}^{(k)}$ ,  $\mathbf{W}^{(v)}$  が注意の規準を定める（学習後は規準固定）

- 出力  $\mathbf{y}_n$  は入力  $\mathbf{v}_n$  の線形和

$$\mathbf{y}_n = \sum_{m=1}^N a_{nm} \mathbf{v}_m$$

- $\mathbf{v}_n$  は、 $(\mathbf{W}^{(o)})^\top$  の列ベクトルの線形和

$$\mathbf{v}_n^\top = \mathbf{x}_n^\top \mathbf{W}^{(o)} \rightarrow \mathbf{v}_n = (\mathbf{W}^{(o)})^\top \mathbf{x}_n$$

つまり

$$\mathbf{v}_n = \sum_{i=1}^D x_{ni} \mathbf{w}_i^{(o)}$$

ただし、 $\mathbf{x}_n = [x_{n1}, x_{n2}, \dots, x_{nD}]^\top$ 、

$$(\mathbf{W}^{(o)})^\top = [\mathbf{w}_1^{(o)}, \mathbf{w}_2^{(o)}, \dots, \mathbf{w}_D^{(o)}].$$



# 注意機構の性質

重要： $\mathbf{W}^{(q)}$ ,  $\mathbf{W}^{(k)}$ ,  $\mathbf{W}^{(v)}$  が注意の規準を定める（学習後は規準固定）

- 出力  $\mathbf{y}_n$  は入力  $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_D$  が張る部分空間に限定される

$$\mathbf{y}_n = \sum_{m=1}^N a_{nm} \mathbf{v}_m$$

- 出力  $\mathbf{y}_n$  は入力  $\mathbf{x}_n$  の線形変換
- 出力  $\mathbf{y}_n$  は学習で得られる  $\mathbf{w}_1^{(o)}, \mathbf{w}_2^{(o)}, \dots, \mathbf{w}_D^{(o)}$  が張る部分空間に限定される。
- 非線形性は注意の強度,  $a_{nm}$ , 経由で導入 ( $\mathbf{x}_n$  に依存・Softmax)

重要： $\mathbf{W}^{(q)}$ ,  $\mathbf{W}^{(k)}$ ,  $\mathbf{W}^{(v)}$  が注意の規準を定める（学習後は規準固定）

- $a_{nm}$  ( $m = 1, 2, \dots, N$ ) は、スケール  $\sqrt{D_k}$  で調整したあとでも、ひとつの要素だけ 1 に近い値となり、その他はゼロに近い値となることが多い。
  - ソフトマックス関数の  $\exp$  が主要因
  - 注意機構は「特定の規準のみ」で注意の大小を決める傾向あり

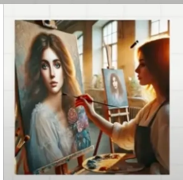
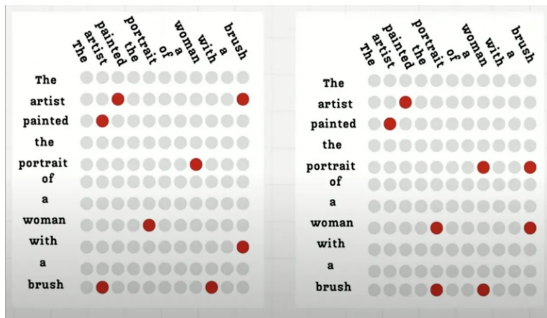
重要： $W^{(q)}$ ,  $W^{(k)}$ ,  $W^{(v)}$  が注意の規準を定める（学習後は規準固定）

例：The artist painted the portrait of a woman with a brush.

上記の文の意味はどっち？

- “a woman with a brush” を描いた
- “a woman” をブラシ（筆）を使って描いた

# Attention Head の注意の規準



case 1 : Painting a woman with a "brush"



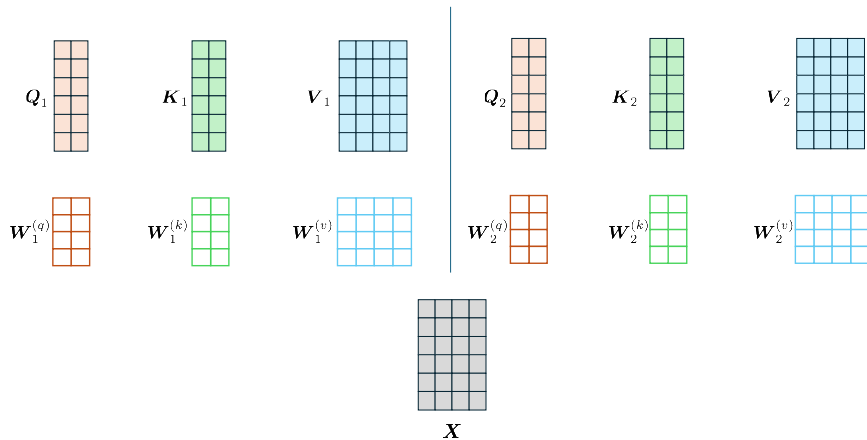
case 2: Painting of a "woman with a brush"

<https://medium.com/@shravankoninti/a-visual-explanation-of-multi-head-attention-6399d86fe51c>

# Attention Head を複数

- 注意の強弱を決める基準はひとつでは足りない
- 時制 / 語彙 / 語順 / ...
- Multi-Head Attention: 複数の基準 ( $\mathbf{W}^{(q)}$ ,  $\mathbf{W}^{(k)}$ ,  $\mathbf{W}^{(v)}$ ) を用意

# マルチヘッド注意機構



# マルチヘッド注意機構

$H$  個の注意機構を用いることにする。 $h$  番目 ( $h = 1, 2, \dots, H$ ) の注意機構は次式でトークンを変換する。

$$\mathbf{H}_h = \text{Attention}(\mathbf{Q}_h, \mathbf{K}_h, \mathbf{V}_h)$$

クエリ・キー・バリューは次式で計算する

$$\mathbf{Q}_h = \mathbf{X}\mathbf{W}_h^{(q)},$$

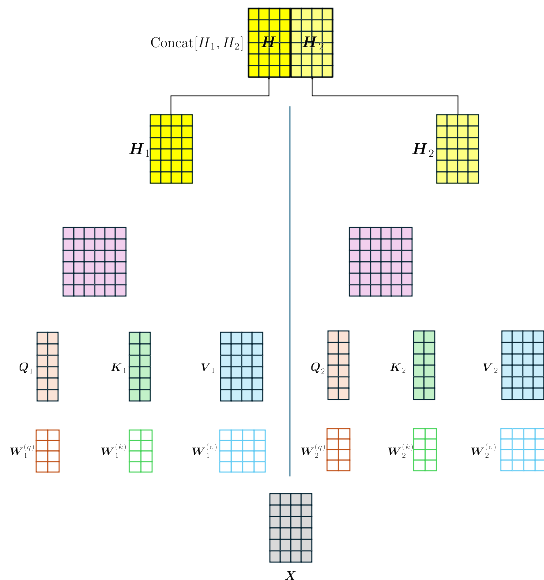
$$\mathbf{K}_h = \mathbf{X}\mathbf{W}_h^{(k)},$$

$$\mathbf{V}_h = \mathbf{X}\mathbf{W}_h^{(v)}.$$

$$\underline{\mathbf{Y}(\mathbf{X}) = \text{Concat}[\mathbf{H}_1, \mathbf{H}_2, \dots, \mathbf{H}_H]\mathbf{W}^{(o)}}$$

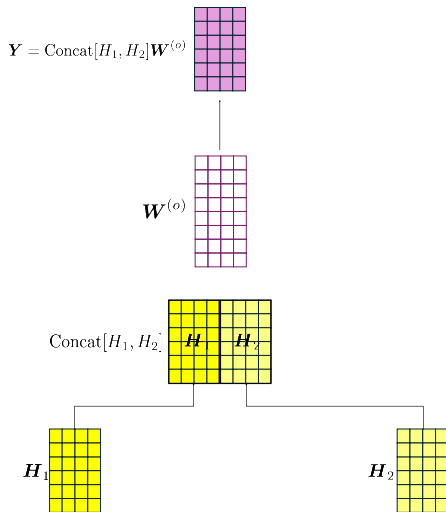
ただし、 $\mathbf{H}_h \in \mathbb{R}^{N \times D_v}$ ,  $\text{Concat}[\mathbf{H}_1, \dots, \mathbf{H}_H] \in \mathbb{R}^{N \times HD_v}$ ,  $\mathbf{W}^{(o)} \in \mathbb{R}^{HD_v \times D}$ ,  $\mathbf{Y}(\mathbf{X}) \in \mathbb{N} \times \mathbb{D}$ 。多くの場合、 $D_v = D/H$

# マルチヘッド注意機構

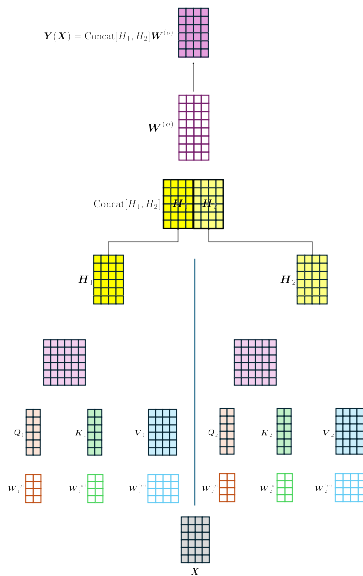




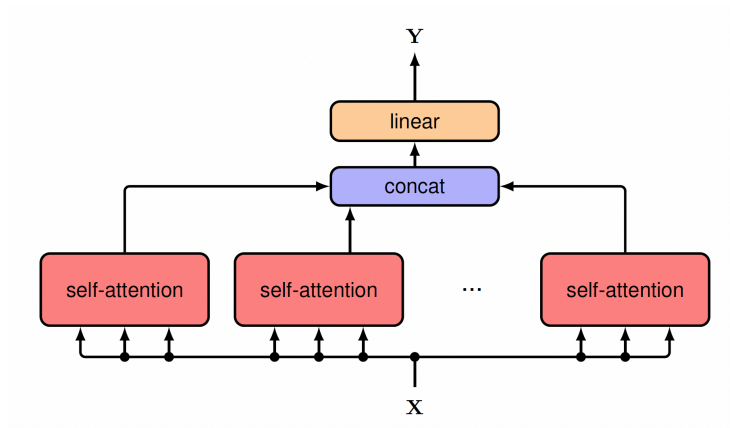
# マルチヘッド注意機構



# マルチヘッド注意機構



# マルチヘッド注意機構



# 準備：レイヤ正規化 (Layer Normalization)

- 誤差逆伝播の計算を安定化させる効果あり  
それぞれの特徴ベクトル  $\mathbf{y}_n \in \mathbb{R}^D$  を「正規化」する。

$$\mathbf{y}_n = [y_{n1}, y_{n2}, \dots, y_{nD}]^\top$$

準備：平均と分散

$$\mu_n = \frac{1}{D} \sum_{i=1}^D y_{ni}$$

$$\sigma_n^2 = \frac{1}{D} \sum_{i=1}^D (y_{ni} - \mu_n)^2$$

正規化

$$\hat{y}_{ni} = \frac{y_{ni} - \mu_n}{\sqrt{\sigma_n + \delta}}$$

記法

$$\text{LayerNorm}[\mathbf{y}_n] = [\hat{y}_{n1}, \hat{y}_{n2}, \dots, \hat{y}_{nD}]^\top$$

学習効率を上げる工夫: レイヤ正規化と、残差表現

マルチヘッドの出力:  $\mathbf{Y}(\mathbf{X})$

$$\mathbf{Y}(\mathbf{X}) = \text{Concat}[\mathbf{H}_1, \mathbf{H}_2, \dots, \mathbf{H}_H] \mathbf{W}^{(o)}$$

元の入力を足して正規化。 $\mathbf{Y}(\mathbf{X})$  は元との差分のみ表現すれば良かった。

$$\mathbf{Z} = \text{LayerNorm}[\mathbf{Y}(\mathbf{X}) + \mathbf{X}]$$

注意機構の出力は  $\mathbf{v}$  が張る線形部分空間に限定。柔軟性を改善するために Multi-Layer Perceptron を導入。例えば、全結合層を 2 つ追加する。

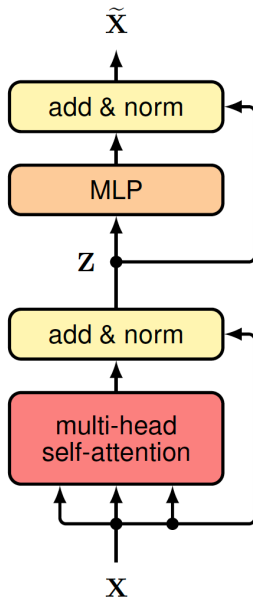
$$\mathbf{Z}' = \text{MLP}[\mathbf{Z}]$$

- 同じ MLP を  $\mathbf{Z}$  の各行に適用する。つまり、各単語のトークンに同じ MLP を適用する。

学習効率を上げる工夫：レイヤ正規化と、残差表現

$$\tilde{\mathbf{X}} = \text{LayerNorm}[\mathbf{Z}' + \mathbf{Z}] = \text{LayerNorm}[\text{MLP}[\mathbf{Z}] + \mathbf{Z}]$$

# トランスフォーマーのアーキテクチャ (符号器・1層分)





# 位置の表現

トークンの入力順序が変わっても出力は変わらない。このままでは、位置・順序の情報が結果に反映されない。

(再掲: 系列の先頭から  $n$  番目のトークンの変換)

$$\mathbf{y}_n = \sum_{m=1}^N a_{nm} \mathbf{x}_m,$$
$$a_{nm} = \frac{\exp\left(\frac{\mathbf{q}_n^\top \mathbf{k}_m}{\sqrt{D_k}}\right)}{\sum_k \exp\left(\frac{\mathbf{q}_n^\top \mathbf{k}_k}{\sqrt{D_k}}\right)}.$$

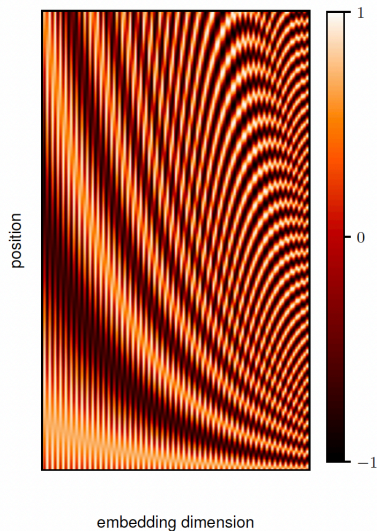
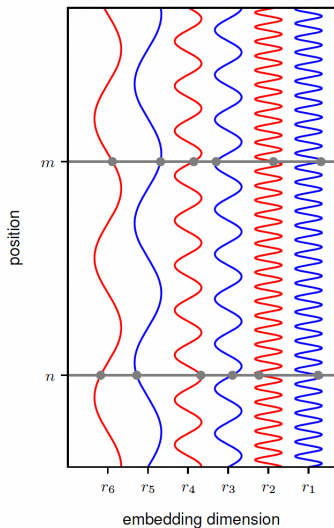
位置の表現  $\mathbf{r}_n = [r_{n1}, r_{n2}, \dots, r_{nD}]^\top$  を生成して、入力トークンに足す

$$\tilde{\mathbf{x}}_n = \mathbf{x}_n + \mathbf{r}_n$$

最大系列長より充分大きい数,  $L$  (例えば 10,000)

$$r_{ni} = \begin{cases} \sin\left(\frac{n}{L^{(i-1)/D}}\right), & \text{if } i \text{ is odd,} \\ \cos\left(\frac{n}{L^{(i-2)/D}}\right), & \text{if } i \text{ is even.} \end{cases}$$

# 位置の表現



# 位置の表現

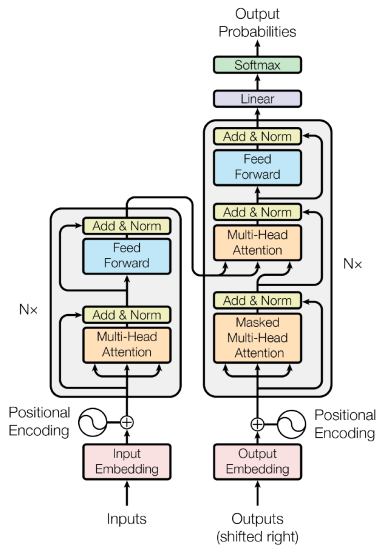
## 位置の表現 $\mathbf{r}$ の性質

- $\|\mathbf{r}_n\| = D/2$
- 位置が一意的に表現される
- 近い位置では似たベクトル、遠い位置では似ていないベクトル
- $\mathbf{r}_n$  と  $\mathbf{r}_{n+k}$  の関係は  $n$  に依らない ( $k$  だけに依存する) 線形変換で表現可能 ( $\phi = k/L^{j/D}$ )

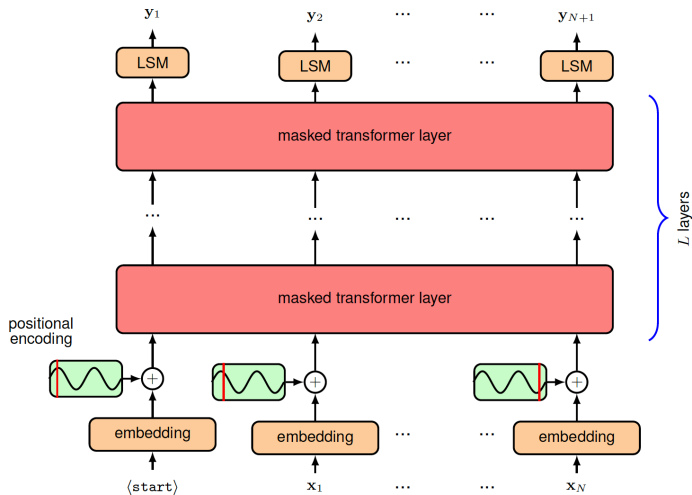
$$\begin{aligned}\sin\left(\frac{n+k}{L^{j/D}}\right) &= \sin\left(\frac{n}{L^{j/D}} + \phi\right) \\ &= \cos(\phi) \sin\left(\frac{n}{L^{j/D}}\right) + \sin(\phi) \cos\left(\frac{n}{L^{j/D}}\right)\end{aligned}$$

- 系列内での相対位置が重要なときには、ニューラルネットワークが ( $\mathbf{W}^{(q)}$ ,  $\mathbf{W}^{(k)}$ ) で学習可能

# トランスフォーマのアーキテクチャ



# 予告：GPT のアーキテクチャ



LSM: Linear SoftMax = 線形多クラス識別器