

Three green apples are arranged in a cluster. One apple is in the foreground, slightly to the right, and is the most prominent. Behind it and to the left are two other apples. The apples are bright green with some lighter green highlights, suggesting a glossy surface. The background is a plain, light color.

人工知能

IE229 – ARTIFICIAL INTELLIGENCE

第14回講義: 自然言語処理 (3)

Lecture 14: Natural Language Processing (3)

二宮 崇 (Takashi Ninomiya)

愛媛大学 (Ehime University)

ninomiya@cs.ehime-u.ac.jp

今回の講義内容

- **自然言語処理 (Natural Language Processing)**
 - ニューラル機械翻訳 (Neural Machine Translation)
 - アテンション付きエンコーダー・デコーダーモデル
 - Transformer



アテンション付きエンコー
ダー・デコーダーモデル
(ENCODER-DECODER MODEL
WITH ATTENTION)

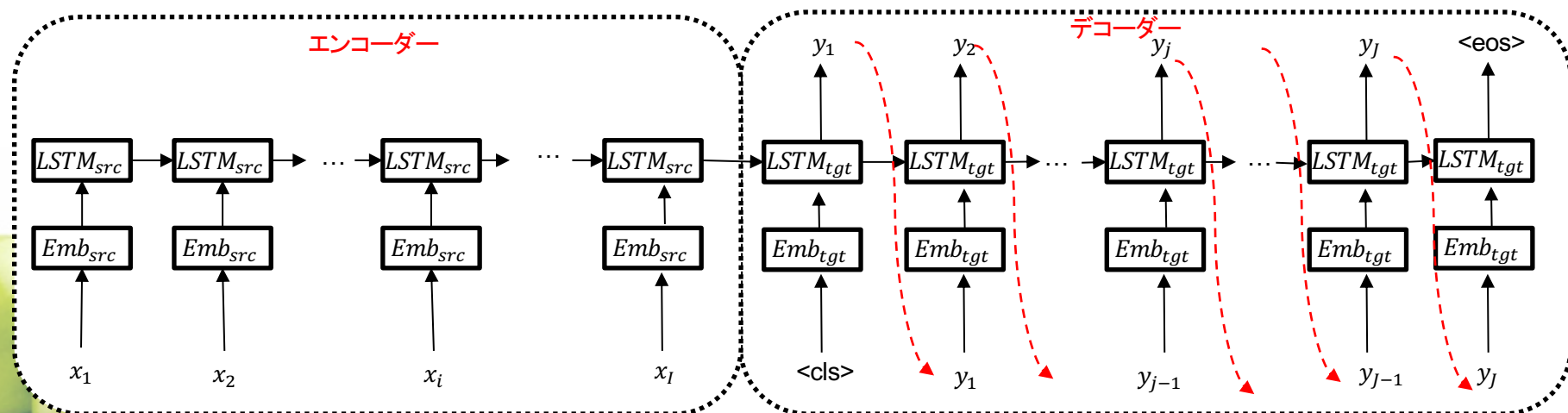


LSTMに基づく エンコーダー・デコーダーモデル

- LSTMに基づくエンコーダー・デコーダーモデル

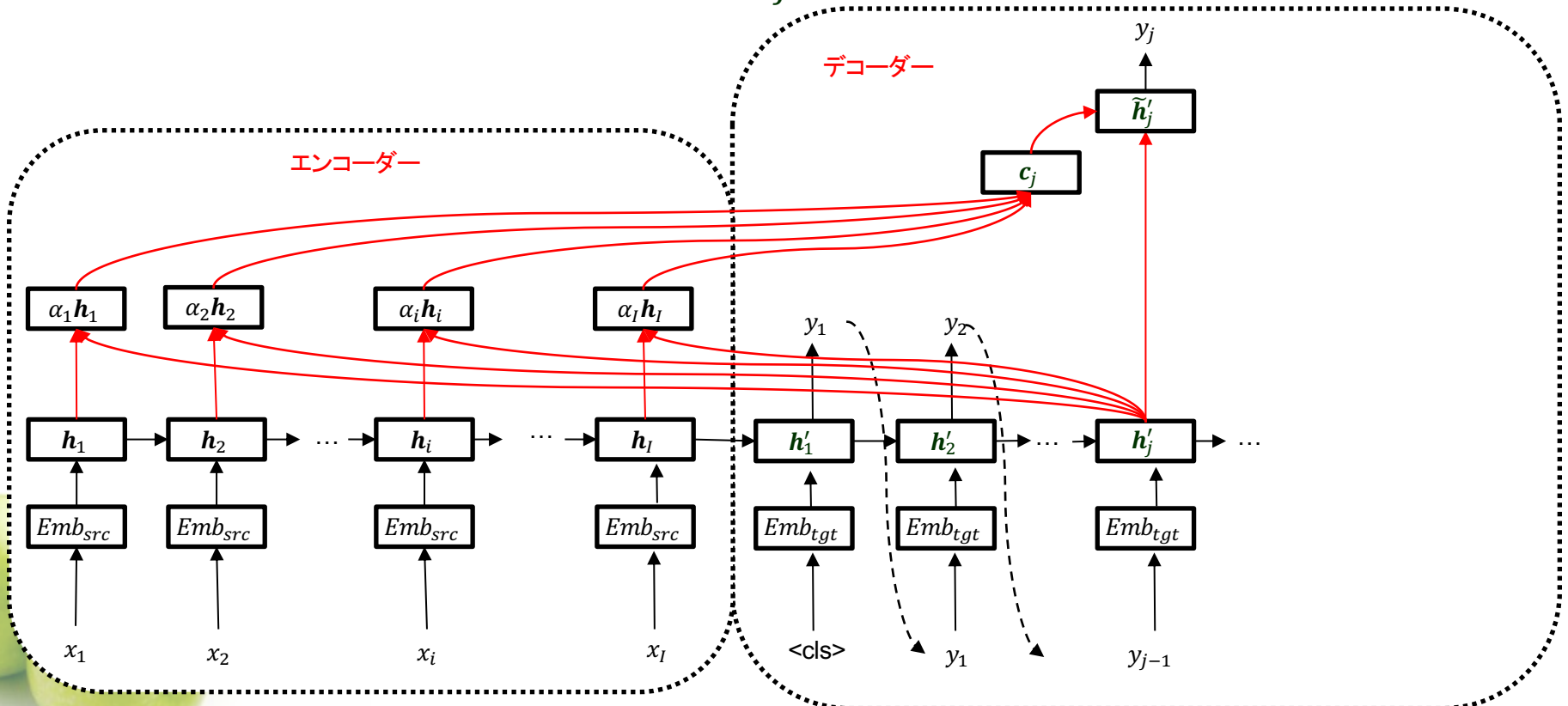
- 入力文: $x_1, x_2, \dots, x_i, \dots, x_I$ 、出力文: $y_1, y_2, \dots, y_j, \dots, y_J$
- エンコーダー: $\mathbf{h}_i = \text{LSTM}_{src}(\text{Emb}_{src}(x_i), \mathbf{h}_{i-1})$
- デコーダー: $\mathbf{h}_j = \text{LSTM}_{tgt}(\text{Emb}_{tgt}(y_{j-1}), \mathbf{h}_{j-1})$
 $y_j = \text{softmax}(W_o \mathbf{h}_j + \mathbf{b}_o)$

- しかし、LSTMを使っても、エンコーダーとデコーダーの間でかなりの距離があるので情報がうまく伝播されない



アテンション付き エンコーダー・デコーダーモデル

- デコーダーの出力の際にエンコーダーの内部状態を直接参照することで長距離の問題を解決する
- デコーダーの内部状態(h'_j)に対するエンコーダーの内部状態(h_i)の重み(α_i)を求め、 h_i の加重平均(c_j ; コンテキストベクトル)を求める



アテンション付き エンコーダー・デコーダーモデル

- アテンション付きエンコーダー・デコーダーモデルの形式化

- 入力文: $x_1, \dots, x_i, \dots, x_I$

- 出力文: $y_1, \dots, y_j, \dots, y_J$

- エンコーダー: $\mathbf{h}_i = LSTM_{src}(Emb_{src}(x_i), \mathbf{h}_{i-1})$ ($\mathbf{h}_i \in \mathbb{R}^d$)

- デコーダー: $\mathbf{h}'_j = LSTM_{tgt}(Emb_{tgt}(y_{j-1}), \mathbf{h}'_{j-1})$ ($\mathbf{h}'_j \in \mathbb{R}^d$)

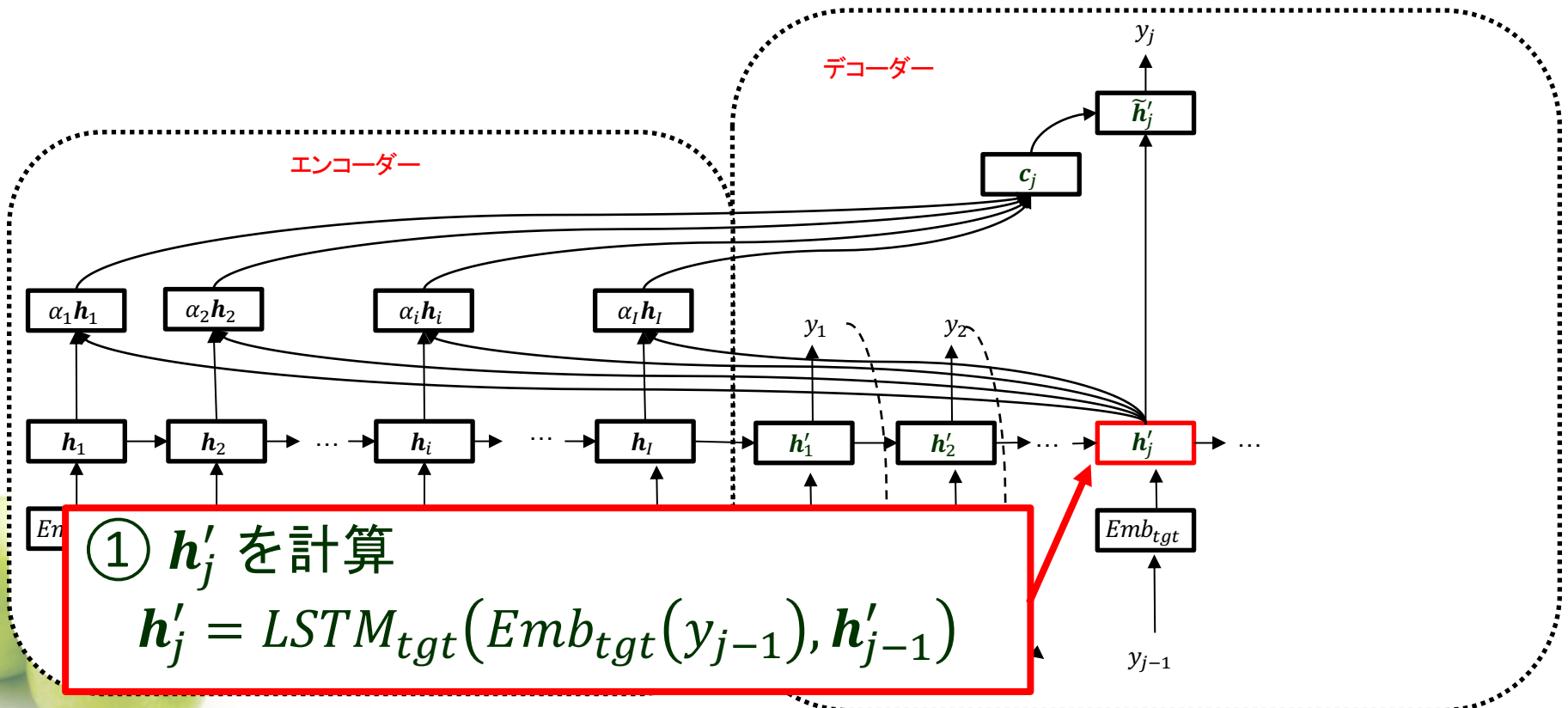
アテンションの重み $\alpha_i = \frac{\exp(\mathbf{h}_i^T \mathbf{h}'_j)}{\sum_{i'=1}^I \exp(\mathbf{h}_{i'}^T \mathbf{h}'_j)}$ ($\alpha_i \in \mathbb{R}$)

コンテキストベクトル $\mathbf{c}_j = \sum_{i=1}^I \alpha_i \mathbf{h}_i$ ($\mathbf{c}_i \in \mathbb{R}^d$)

$$\tilde{\mathbf{h}}'_j = \tanh(W_c[\mathbf{c}_j; \mathbf{h}'_j] + \mathbf{b}_c)$$
 ($W_c \in \mathbb{R}^{d \times 2d}$)
$$y_j = \text{softmax}(W_o \tilde{\mathbf{h}}'_j + \mathbf{b}_o)$$
 ($W_o \in \mathbb{R}^{V \times d}$)

;は連結、 d は内部状態の次元数、 V は語彙数

アテンション付き エンコーダー・デコーダーモデル

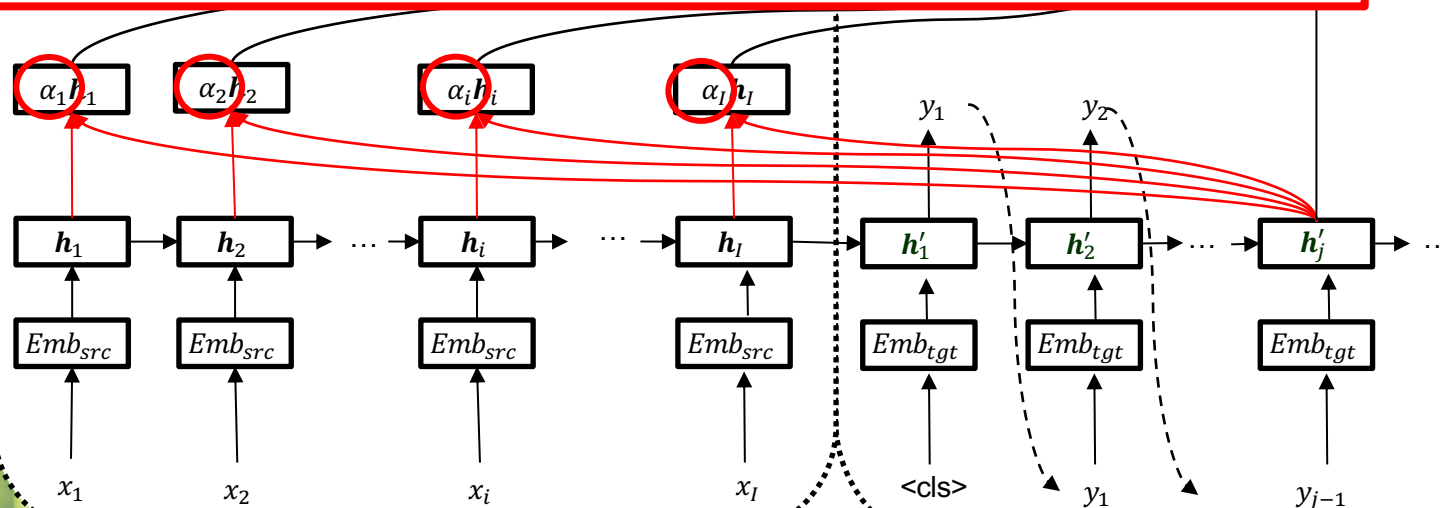


アテンション付き エンコーダー・デコーダーモデル

② h_i と h'_j から α_i を計算($1 \leq i \leq I$)

$$\alpha_i = \frac{\exp(\mathbf{h}_i^T \mathbf{h}'_j)}{\sum_{i'=1}^I \exp(\mathbf{h}_{i'}^T \mathbf{h}'_j)}$$

- α_i は h_i と h'_j の内積にsoftmaxをかけたもの
- $0 \leq \alpha_i \leq 1, \sum_{i=1}^I \alpha_i = 1$ となることに注意

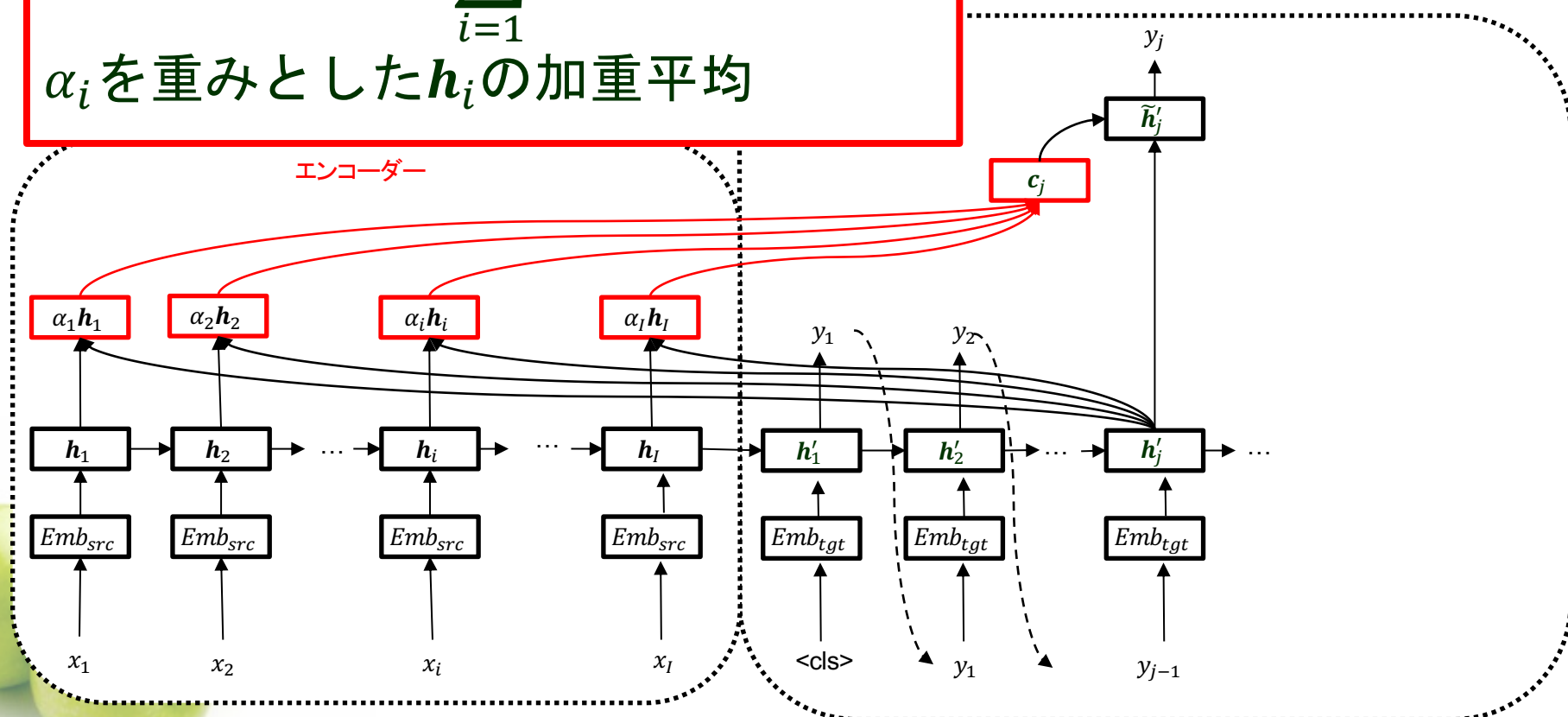


アテンション付き エンコーダー・デコーダーモデル

③ コンテキストベクトル c_j を計算

$$c_j = \sum_{i=1}^I \alpha_i h_i$$

α_i を重みとした h_i の加重平均

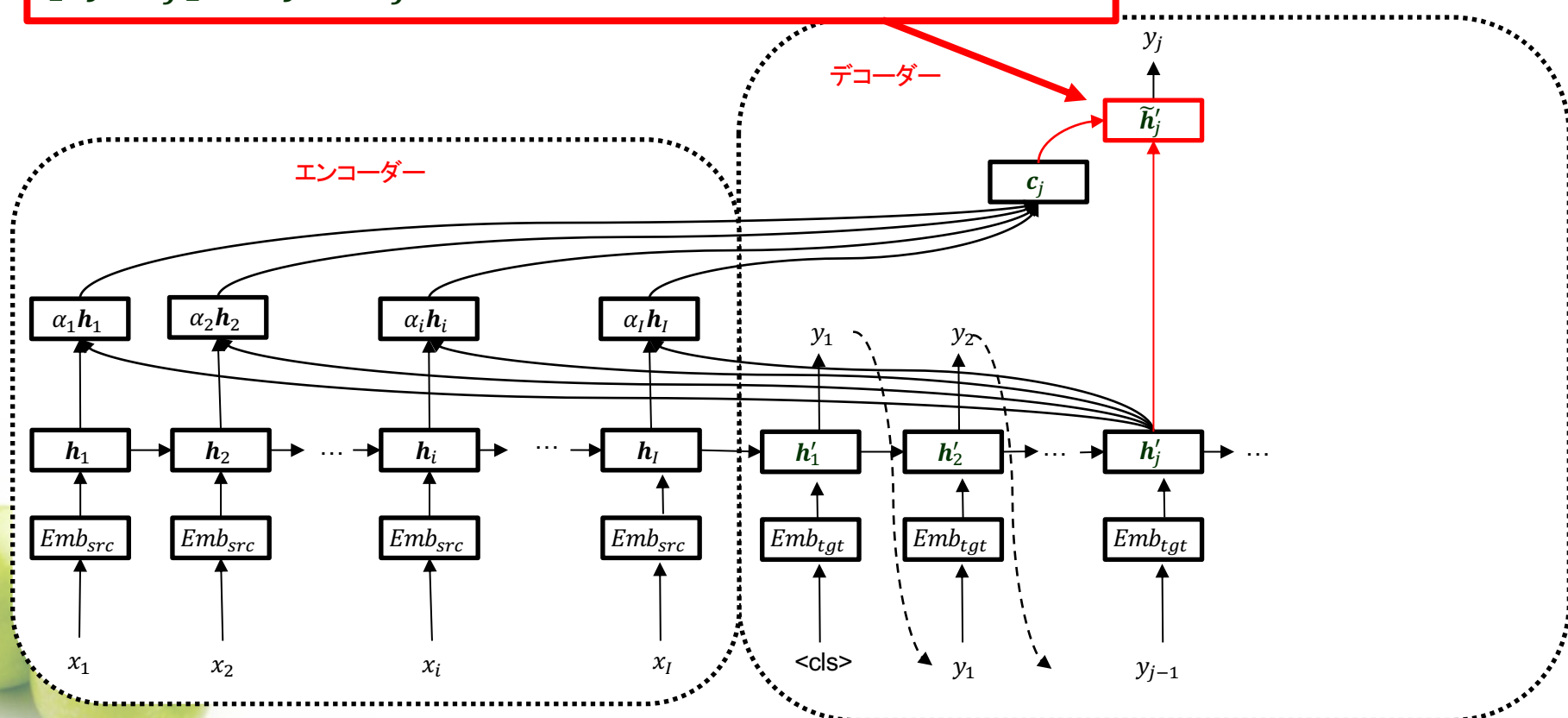


アテンション付き エンコーダー・デコーダーモデル

④ 新しい内部状態 \tilde{h}'_j を計算

$$\tilde{h}'_j = \tanh(W_c[c_j; h'_j] + b_c)$$

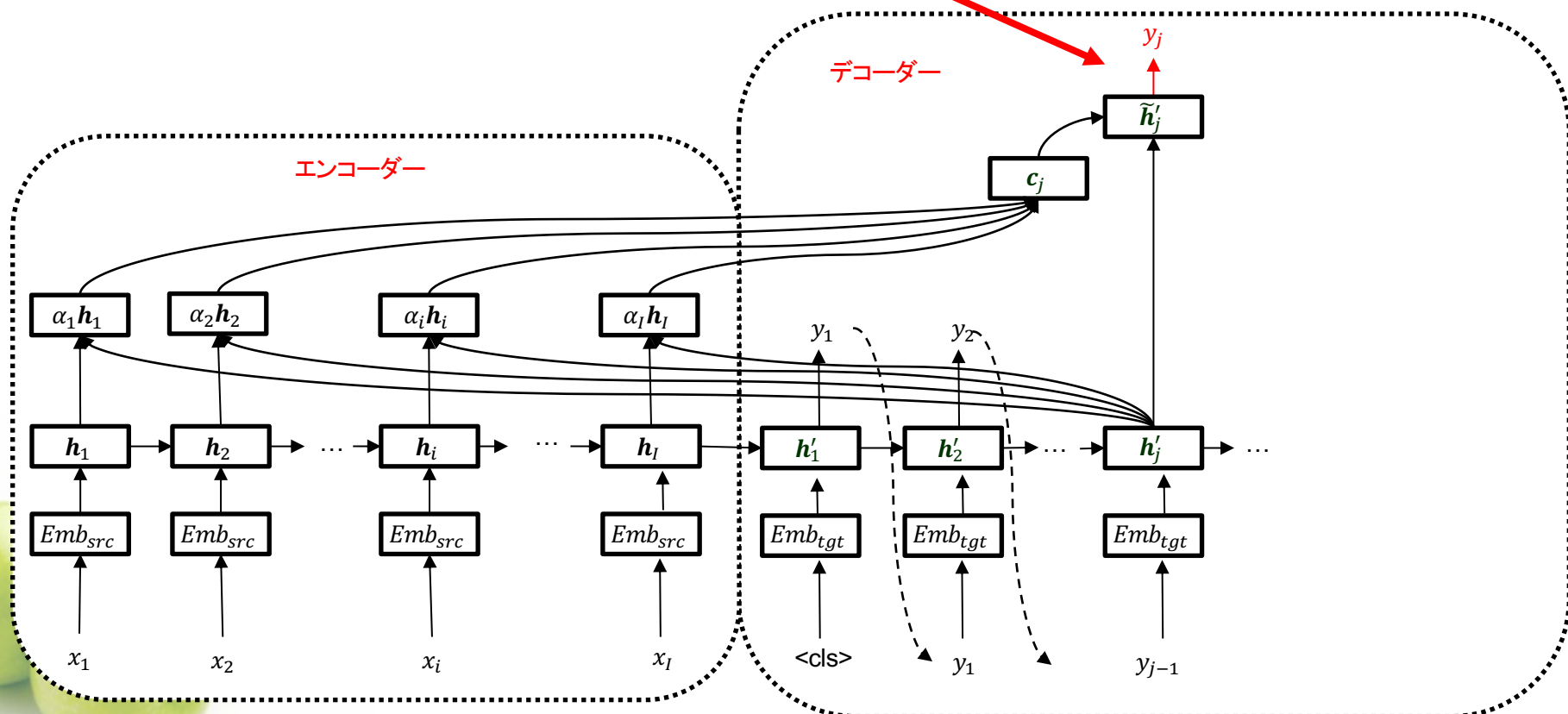
$[c_j; h'_j]$ は c_j と h'_j を連結したベクトル



アテンション付き エンコーダー・デコーダーモデル

⑤出力を計算

$$y_j = \text{softmax}(W_o \tilde{h}'_j + b_o)$$



アテンション付き エンコーダー・デコーダーモデル

● まとめ

- エンコーダーの内部状態に対する重み付け

$$\alpha_i = \frac{\exp(\mathbf{h}_i^T \mathbf{h}'_j)}{\sum_{i'=1}^I \exp(\mathbf{h}_{i'}^T \mathbf{h}'_j)}$$

- 基本的に内積で計算されていることに注意
- $0 \leq \alpha_i \leq 1, \sum_{i=1}^I \alpha_i = 1$ となることに注意
- コンテキストベクトルは、エンコーダーの内部状態の加重平均

$$\mathbf{c}_j = \sum_{i=1}^I \alpha_i \mathbf{h}_i$$

- 新しい内部状態はコンテキストベクトルと内部状態を連結したベクトルに対する線形変換で得られる

$$\tilde{\mathbf{h}}'_j = \tanh(W_c[\mathbf{c}_j; \mathbf{h}'_j] + \mathbf{b}_c)$$



アテンション付き エンコーダー・デコーダーモデル

- 重み付けのバリエーション

$$\alpha_i = \frac{\exp(\text{score}(\mathbf{h}_i, \mathbf{h}'_j))}{\sum_{i'=1}^I \exp(\text{score}(\mathbf{h}_{i'}, \mathbf{h}'_j))}$$

- Global Attention (dot)

$$\text{score}(\mathbf{h}_i, \mathbf{h}'_j) = \mathbf{h}_i^T \mathbf{h}'_j$$

ベクトル間の関係を内積で表現

- 先程説明したものと同一通常のアテンション

- Global Attention (general)

$$\text{score}(\mathbf{h}_i, \mathbf{h}'_j) = \mathbf{h}_i^T W_a \mathbf{h}'_j$$

ベクトル間の関係を重み行列 W_a で学習

- Global Attention (concat)

$$\text{score}(\mathbf{h}_i, \mathbf{h}'_j) = \mathbf{v}^T \tanh(W_a [\mathbf{h}_i; \mathbf{h}'_j])$$

2つのベクトル間の重み付けを重み行列 W_a と重みベクトル \mathbf{v} で学習



TRANSFORMER



Transformer

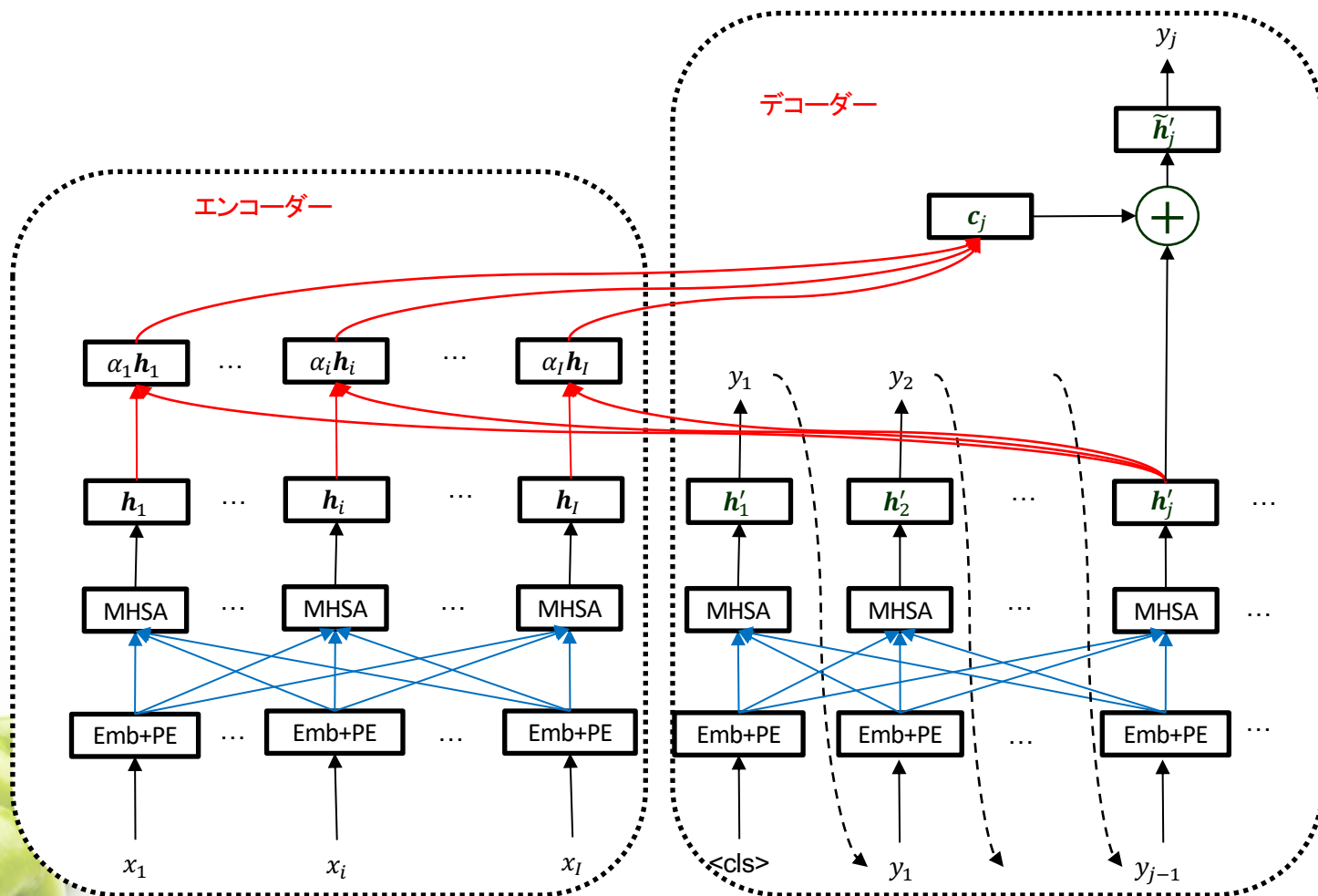
- **Transformerモデル**

- 現在のニューラル機械翻訳のベースラインモデル
- RNNでもCNNでもない新しいタイプのニューラル機械翻訳モデル
- アテンションの計算だけで翻訳を行うことが特徴
- キーテクノロジー
 - 自己アテンション (Self Attention)
 - 位置エンコーディング (Positional Encoding)
 - マルチヘッドアテンション (Multi-head Attention)



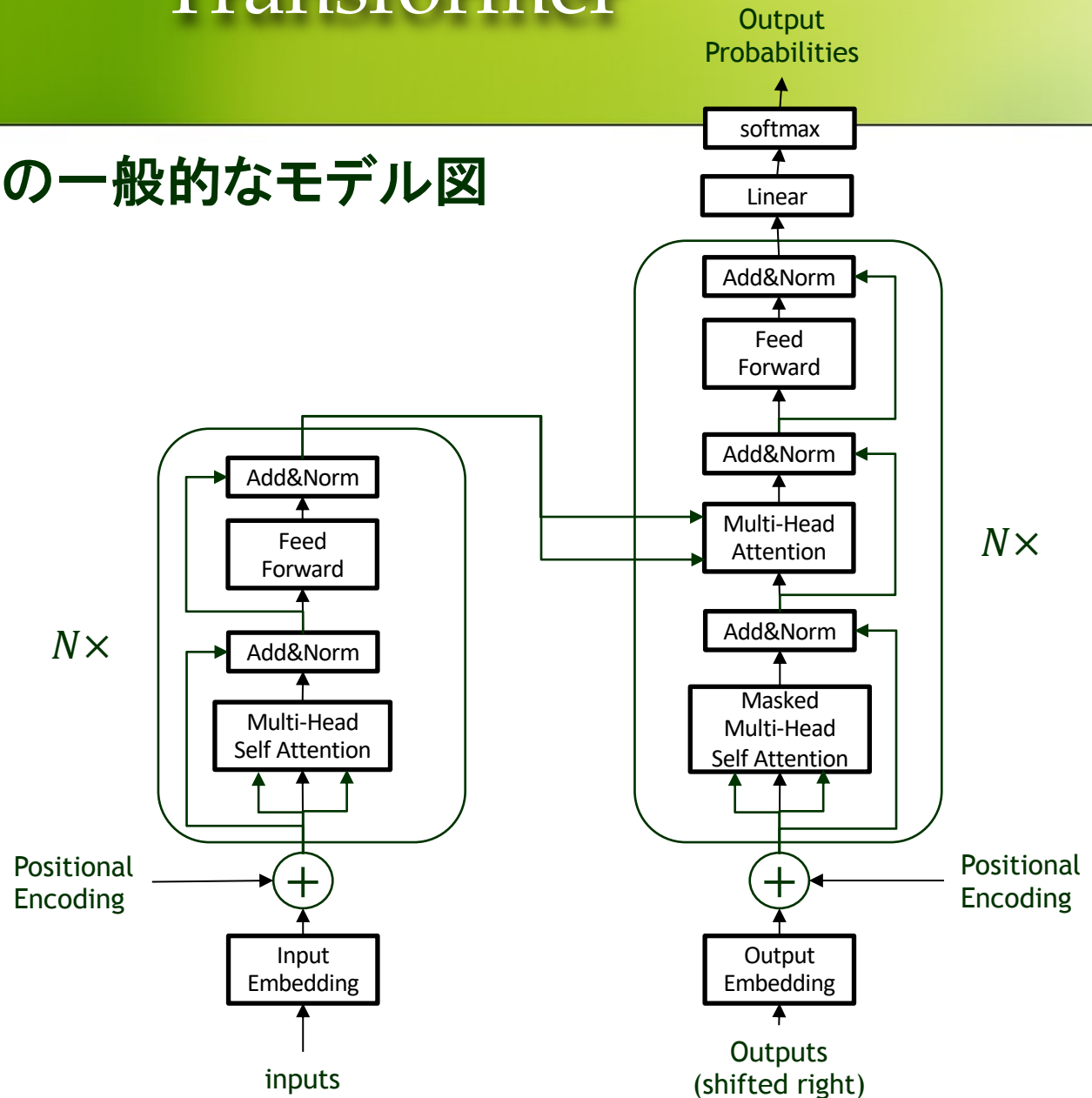
Transformer

- Transformerのモデル概要図



Transformer

- Transformerの一般的なモデル図



Transformer

● アテンション

- d : 内部状態の次元数
- Query, Key, Valueの入力があると考えてアテンションを捉え直す
 - **Query(Q)**: 計算対象の内部状態 (デコーダーの内部状態)
 - **Key(K)**: アテンションの対象 (エンコーダーの内部状態)
 - **Value(V)**: アテンション先の内部状態 (エンコーダーの内部状態)
- **Q 、 K 、 V は(文長 $\times d$)の行列を表し、文全体のアテンションをまとめて計算**
- Q と K から重み(A)を計算して、重みと V の加重平均を計算。 A は(文長 \times 文長)の行列

$$A = \text{softmax} \left(\frac{QK^T}{\sqrt{d}} \right)$$

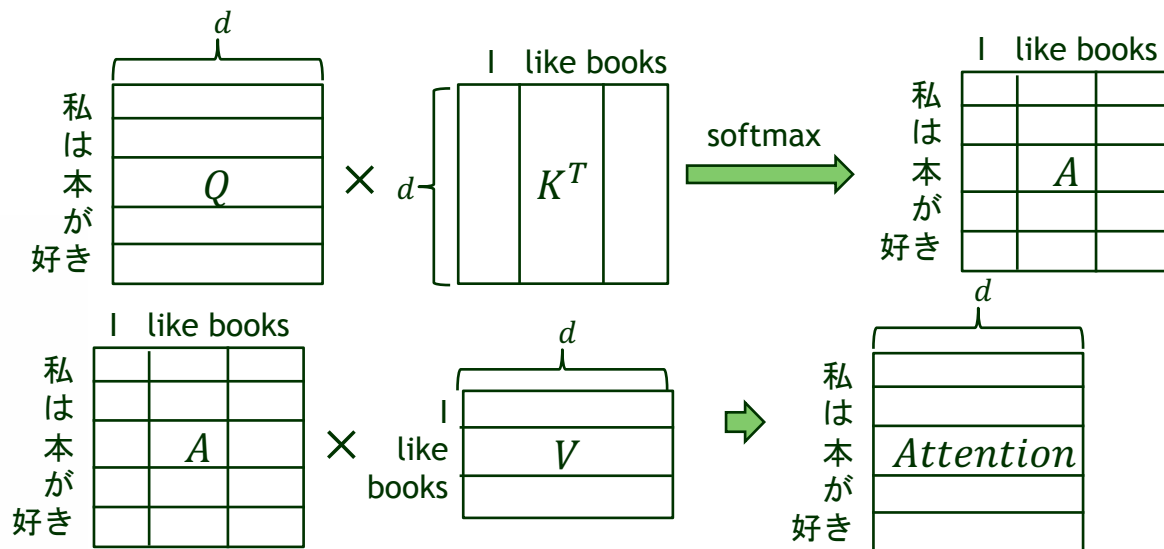
$$\text{Attention}(Q, K, V) = AV$$



Transformer: アテンションの例

- **言語間アテンション**(Encoder-Decoder Attention, Cross-Lingual Attention)

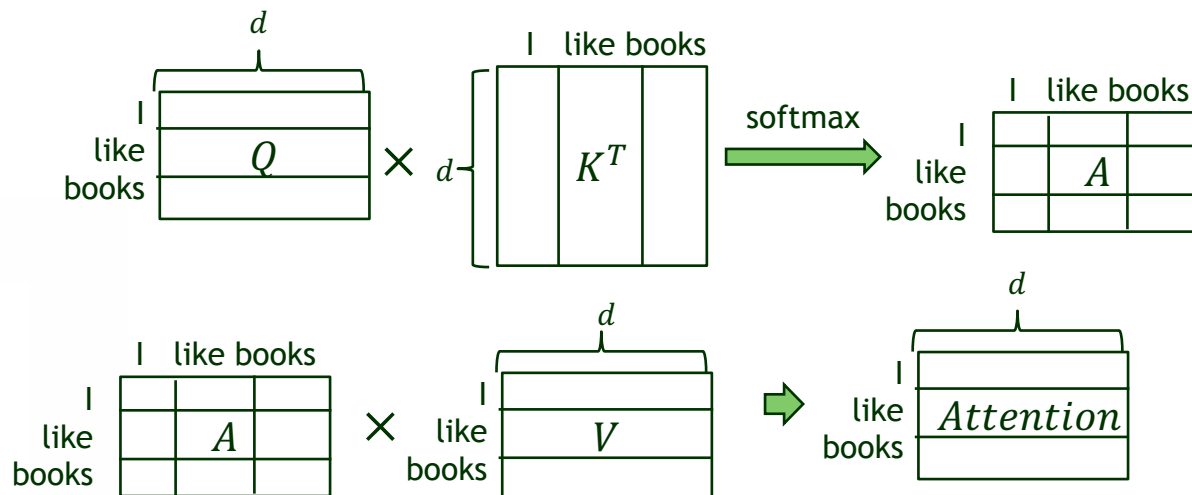
- エンコーダー(“I like books”)→デコーダー(“私は本が好き”)
 - Q : デコーダーの内部状態
 - K : エンコーダーの内部状態
 - V : エンコーダーの内部状態(K と V は同じ)
 - A : QK^T を計算し、各行ごとにsoftmaxの計算



Transformer: 自己アテンション

- 自己アテンション(Self Attention)

- エンコーダーの場合(“I like books”)
 - Q : エンコーダーの内部状態
 - K : エンコーダーの内部状態
 - V : エンコーダーの内部状態(Q と K と V は同じ)
 - A : QK^T を計算し、各行ごとにsoftmaxの計算



Transformer

- マルチヘッドアテンション (Multi-head Attention)

- 内部状態を表すベクトル(\mathbf{h})をヘッド数(H)に分割して、それぞれの分割されたベクトルに対してアテンションの計算を行う
- d : 内部状態の次元数
- H : ヘッド数
 - 例: $d = 512, H = 8$ なら、各ヘッドの次元数は $512/8=64$ 次元)
- 線形変換による分割($W_h^Q, W_h^K, W_h^V \in \mathbb{R}^{d \times (d/H)}$)
$$Q_h = QW_h^Q, \quad K_h = KW_h^K, \quad V_h = VW_h^V$$
- 各ヘッドのアテンションの計算

$$Head_h = softmax\left(\frac{Q_h K_h^T}{\sqrt{d/H}}\right) V_h$$

- アテンションの計算($W^o \in \mathbb{R}^{d \times d}$)
$$Attention(Q, K, V) = [Head_1; \dots; Head_H]W^o$$



Transformer

- **Feed Forward層**

- 各単語ごとの全結合層(Fully Connected Layer)2層+ReLU

$$FFN(x) = \max(0, xW_1 + b_1)W_2 + b_2$$



Transformer

- レイヤー正規化 (Layer Normalization)

- ベクトル内の全ての次元の値の平均、分散を求めて正規化すること
- \mathbf{h} : 内部状態 ($\mathbf{h} = (h_1, \dots, h_d)^T$)
- d : 内部状態の次元数
- γ, β : パラメータ

$$\mu = \frac{1}{d} \sum_{i=1}^d h_i$$
$$\sigma = \sqrt{\frac{1}{d} \sum_{i=1}^d (h_i - \mu)^2}$$

$$\text{LayerNorm}(\mathbf{h}) = \gamma \odot \frac{\mathbf{h} - \mu}{\sigma} + \beta$$

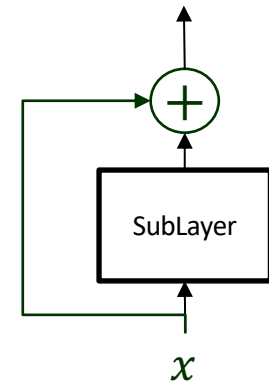


Transformer

- 残差接続(Residual Connection)

- ニューラルネットワークのある部分ネットワークSubLayerがあったとき、次の計算を行うことを残差接続(Residual Connection)という

$$Res(x) = x + SubLayer(x)$$



- Add&Norm層

$$AddNorm(x) = LayerNorm(x + SubLayer(x))$$



Transformer: 位置エンコーディング

- 位置エンコーディング (Positional Encoding)

- 自己アテンションで同じ文内の任意の単語間の関係が計算されているが、単語間の距離が考慮されていない(隣の単語も遠くの単語も同じ様にアテンションの計算が行われる)
- 単語埋め込みの後に、位置情報を埋め込むことでこの問題を解決する

- 位置エンコーディングの式

- pos : 文中の位置
- $2i, 2i + 1$: 位置埋め込みベクトルの次元
- d : 埋め込みベクトルの次元数

$$PE_{(pos, 2i)} = \sin\left(\frac{pos}{10000^{2i/d}}\right)$$

$$PE_{(pos, 2i+1)} = \cos\left(\frac{pos}{10000^{2i/d}}\right)$$



Transformer

Dropoutとラベルスムージング

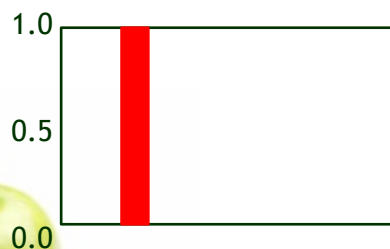
- 各層の後にDropoutを適用
- 出力にラベルスムージング(Label Smoothing)を適用

- 通常の交差エントロピー損失

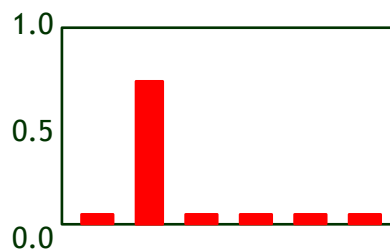
$$L = - \sum_{j=1}^J \sum_{k=1}^K t_{jk} \log P(y_{jk} | x, y_{<j})$$

- ラベルスムージング

$$L = - \sum_{j=1}^J \sum_{k=1}^K \left[(1 - \epsilon) t_{jk} + \epsilon \frac{1}{K} \right] \log P(y_{jk} | x, y_{<j})$$



通常の交差エントロピー



ラベルスムージング

$t_{jk} \in \{0,1\}$ はj番目の出力の
k番目のラベルの正解

ϵ はスムージングの度合いを
調整するハイパーパラメータ



Transformer: デコーダー

- マスク付きマルチヘッドアテンション

- デコーダーの自己アテンションの計算時にマスクをかける
 - 推論時には逐次的に単語が生成されるため、学習時にもマスクをかけて後ろの単語にアテンションがかからないようにする

デコーダーのアテンション行列A

	私	は	本	が	好き
私					
は					
本					
が					
好き					

黒い箇所を自己アテンションの計算対象にしない
(アテンションスコアを0にする)



ニューラル機械翻訳のまとめ

- アテンション付きエンコーダー・デコーダーモデル
- Transformer
 - 自己アテンション
 - マルチヘッドアテンション
 - レイヤー正規化
 - 残差接続
 - 位置エンコーディング
 - マスク付き自己アテンション

