

Deep Learning基礎講座

Natural Language Processing

Day 2

講義シラバス

Day 1: 前処理、Word2Vec

Day 2: Encoder-Decoderモデル

Day 3: Transformerモデル

NLPのカテゴリ例

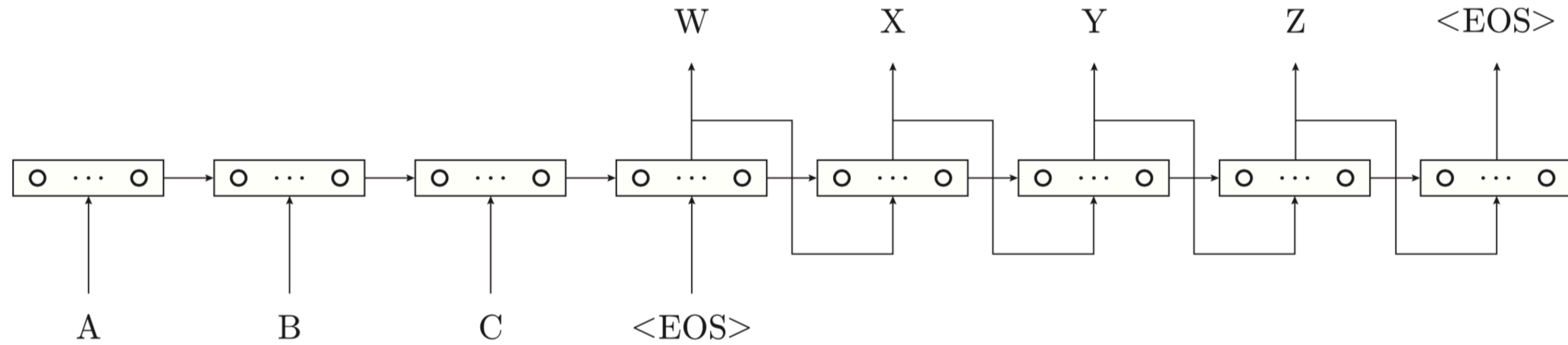
基礎技術： 品詞タグ付け、単語分割、構文解析、固有表現抽出

応用技術： 機械翻訳、文書分類、Q&A、対話

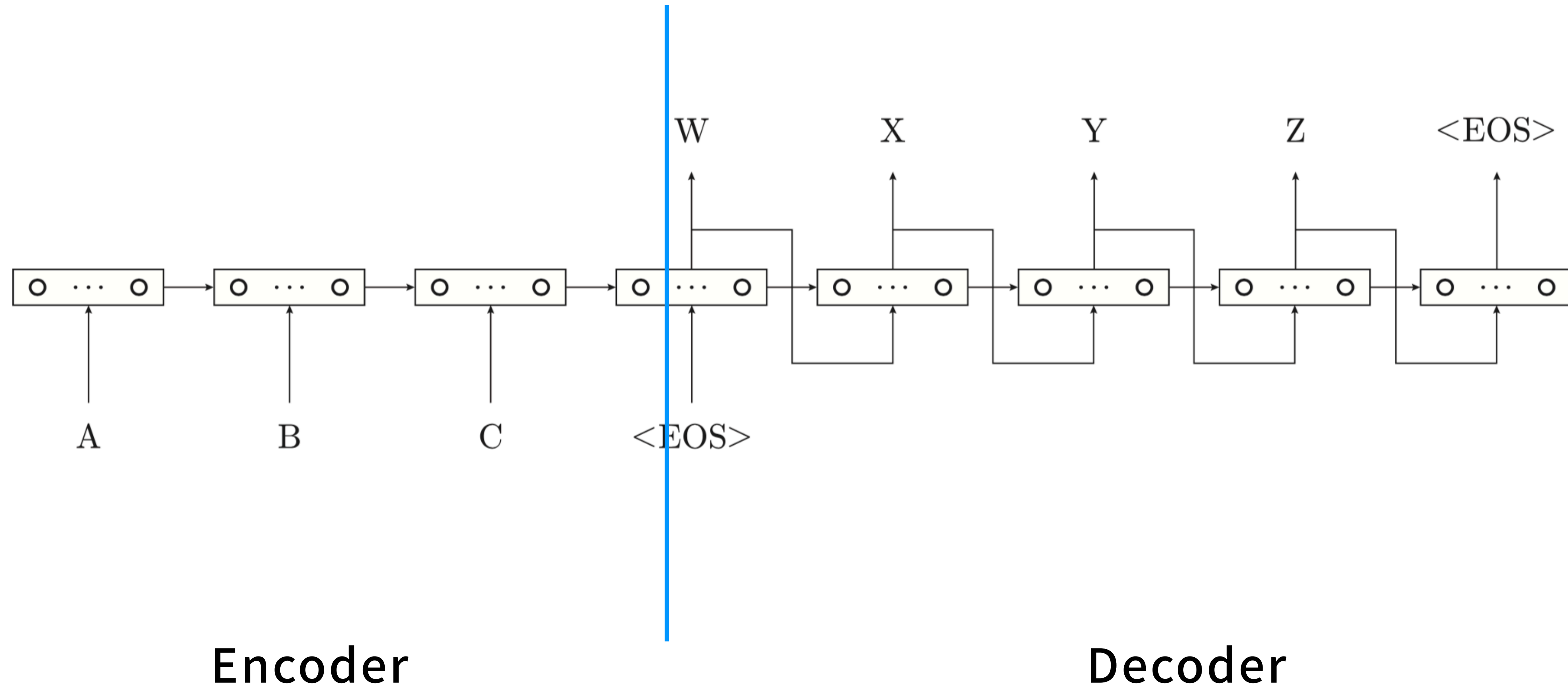


入力も出力も時系列の形式

Encoder-Decoder



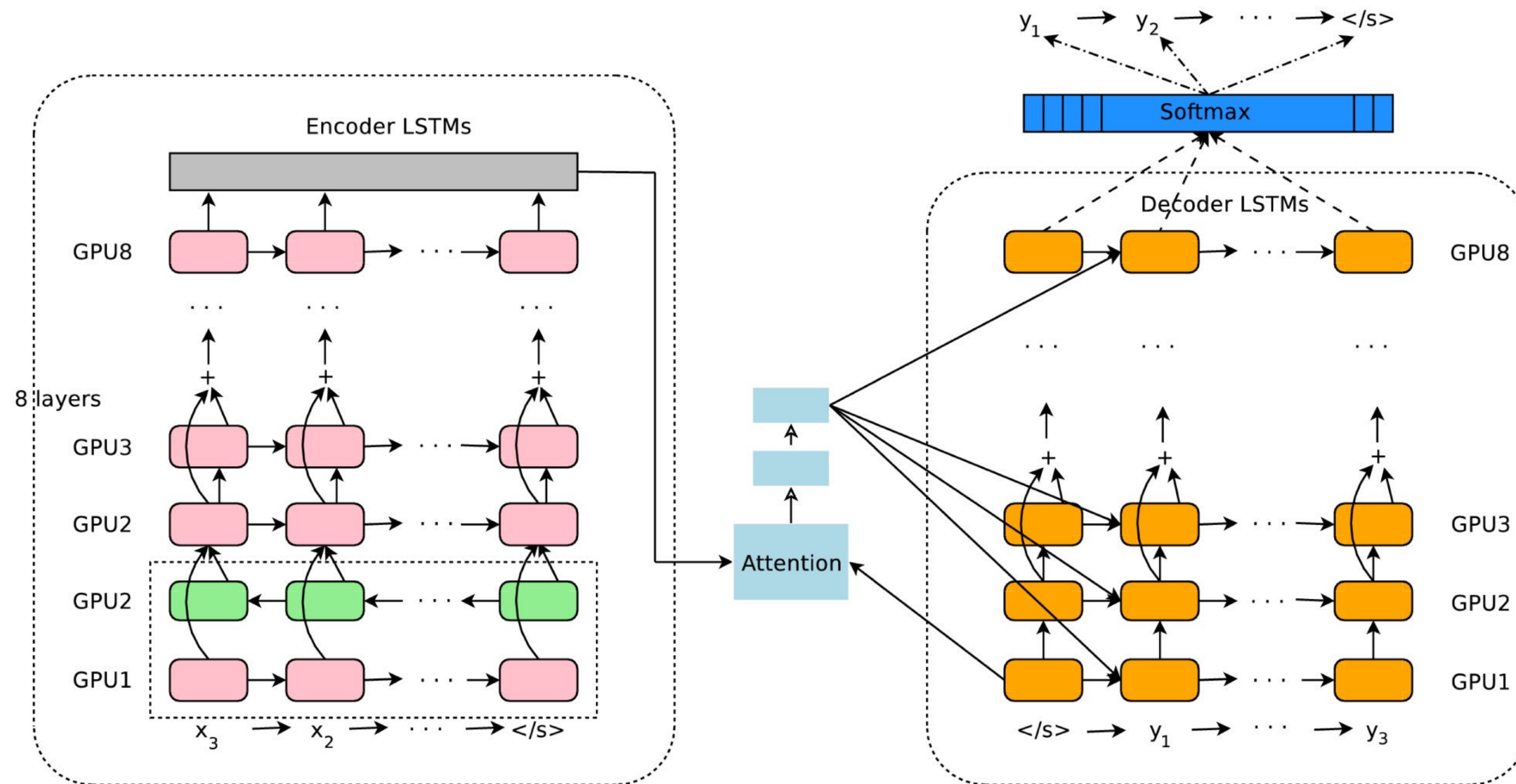
Encoder-Decoder



Encoder-Decoderは2つのネットワークで構成されるモデル

- Encoder: 入力の時系列用のネットワーク
- Decoder: 出力の時系列用のネットワーク

cf. Google翻訳



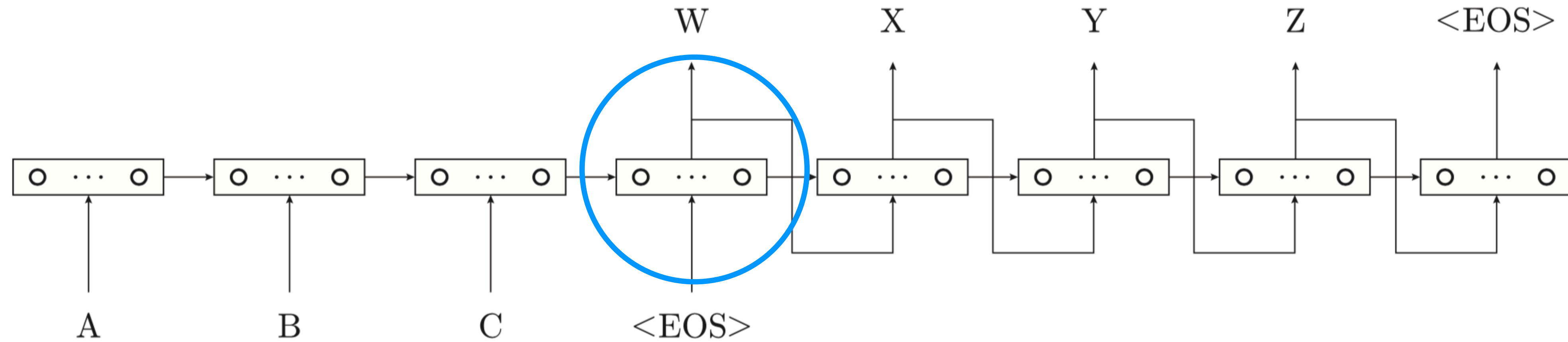
注意：

Encoder, Decoder はRNNであるとは限らない

- 画像キャプションではEncoderはCNN
- 最近ではそもそも翻訳タスクにRNNを用いなくなっている
(詳しくはDay3にて)

Encoder-Decoder の概観

- Encoderの実装はネットワークそのままのケースがほとんど
- 逆に、Decoderではいくつか気をつけるべき点がある



EncoderとDecoderは別ネットワークだが、モデル全体としてはひとつなぎ



Encoderの（隠れ層の）状態は、Decoderの初期値として設定する必要がある

Decoderでは学習の仕方も2通りある

1. Decoderの各時刻の出力を次の入力として学習
2. Decoderへの入力は教師データをそのまま使って学習
(Teacher Forcing)

1. Decoderの各時刻の出力を次の入力として学習
2. Decoderへの入力は教師データをそのまま使って学習



- 1 だと誤差がどんどん大きくなり学習が不安定になる
- 2 だと学習は安定するが、評価時とはデータ分布が異なる

→ **Teacher Forcing** を確率的に取り入れる
Scheduled Sampling を用いる

生成文の評価について：BLEUスコア

- Bilingual Evaluation Understudy（BLUEではなくBLEU）
- 単なる予測誤差ではなく「どれくらい自然な翻訳ができているか」を評価

$$\text{BLEU} = \text{BP} \cdot \exp \left(\sum_{n=1}^N \frac{1}{N} \log p_n \right)$$

$$\text{BLEU} = \underline{\text{BP}} \cdot \exp \left(\sum_{n=1}^N \frac{1}{N} \log p_n \right)$$

Brevity Penalty

$$\text{BP} = \begin{cases} 1 & \text{if } c \geq r \\ \exp(1 - r/c) & \text{if } c < r \end{cases}$$

c : 生成文の長さ (candidate)

r : 正解文の長さ (reference)

→ 短かすぎる文にはペナルティを与える

$$\text{BLEU} = \text{BP} \cdot \exp \left(\sum_{n=1}^N \frac{1}{N} \log p_n \right)$$

Modified n-gram precision

- Reference中の単語は一回しか使えないという制約下でprecisionを計算

例： Candidate: the the the the the the the

Reference: The cat is on the mat

→ Precision = 1, Modified = 1/7

(実際はこれをn-gramに対して計算する)

生成文の評価について：Greedy Search と Beam Search

Greedy:

各時刻で最も確率の高い単語を正解として採用し、次のステップでの入力として使う

Beam:

各時刻において一定の数のそれまでのスコアが高い文を保持しながら選択を行う

(演習で実装あり)

演習

(Define-by-RunのほうがNLPはかなり実装しやすい！)

Attentionについて

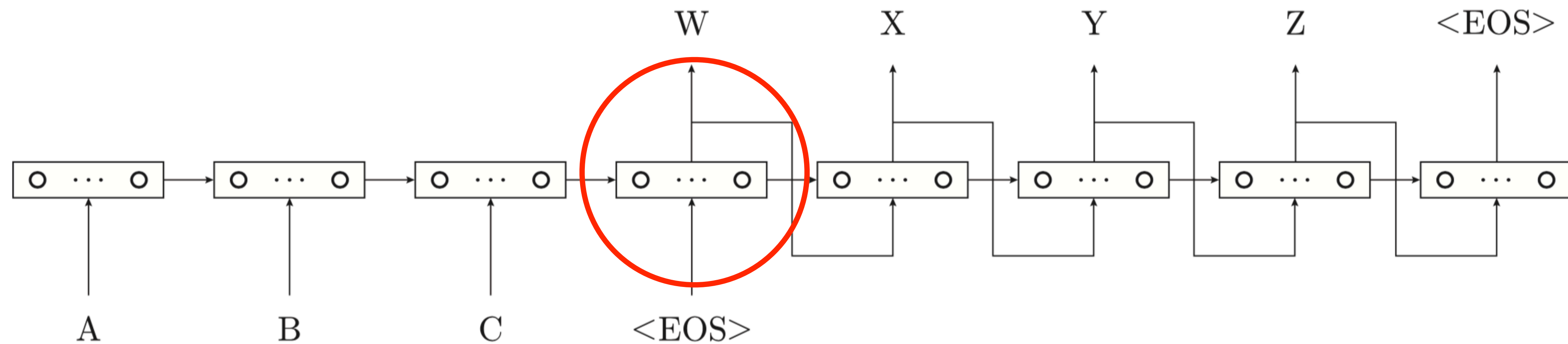
**Encoder-Decoder は一定性能を誇るモデルだが、
よくよく考えると、不要な処理を行っている点がある**

入力系列が与えられたときの出力系列が得られる確率：

$$p(\mathbf{y}(1), \dots, \mathbf{y}(T') \mid \mathbf{x}(1), \dots, \mathbf{x}(T)) = \prod_{t=1}^{T'} p(\mathbf{y}(t) \mid \mathbf{y}(1), \dots, \mathbf{y}(t-1), \mathbf{c})$$

入力系列が与えられたときの出力系列が得られる確率：

$$p(\mathbf{y}(1), \dots, \mathbf{y}(T') \mid \mathbf{x}(1), \dots, \mathbf{x}(T)) = \prod_{t=1}^{T'} p(\mathbf{y}(t) \mid \mathbf{y}(1), \dots, \mathbf{y}(t-1), \mathbf{c})$$



入力を持つ「文脈」情報は、EncoderとDecoderの境界部分に集約されてしまう

- 時系列情報が、結局は固定長のベクトル \mathbf{c} として表現されることになる

しかし、本来は各時刻によって過去のどの時刻を重視すべきなのは異なるはず

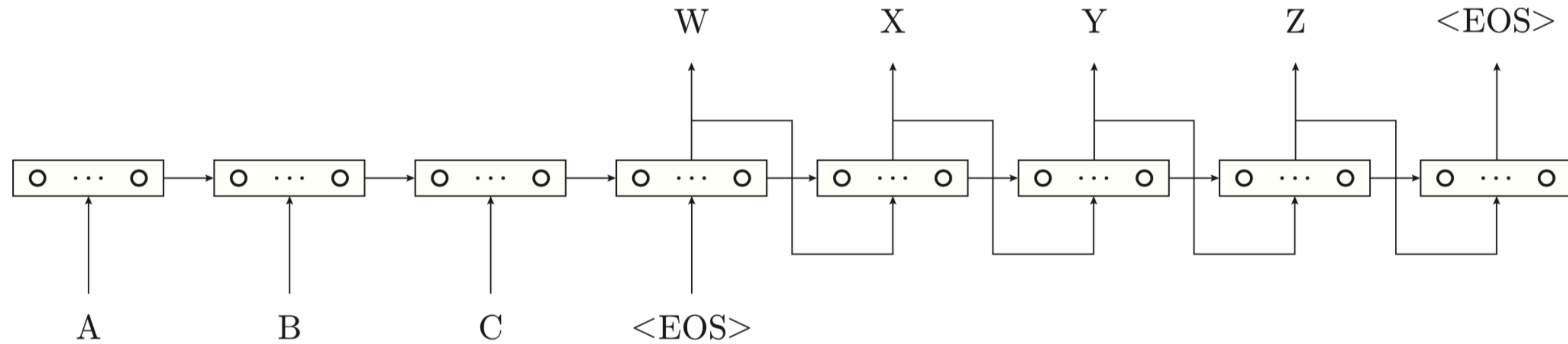


入力系列がもつ情報を1つのベクトルに集約する必然性はない

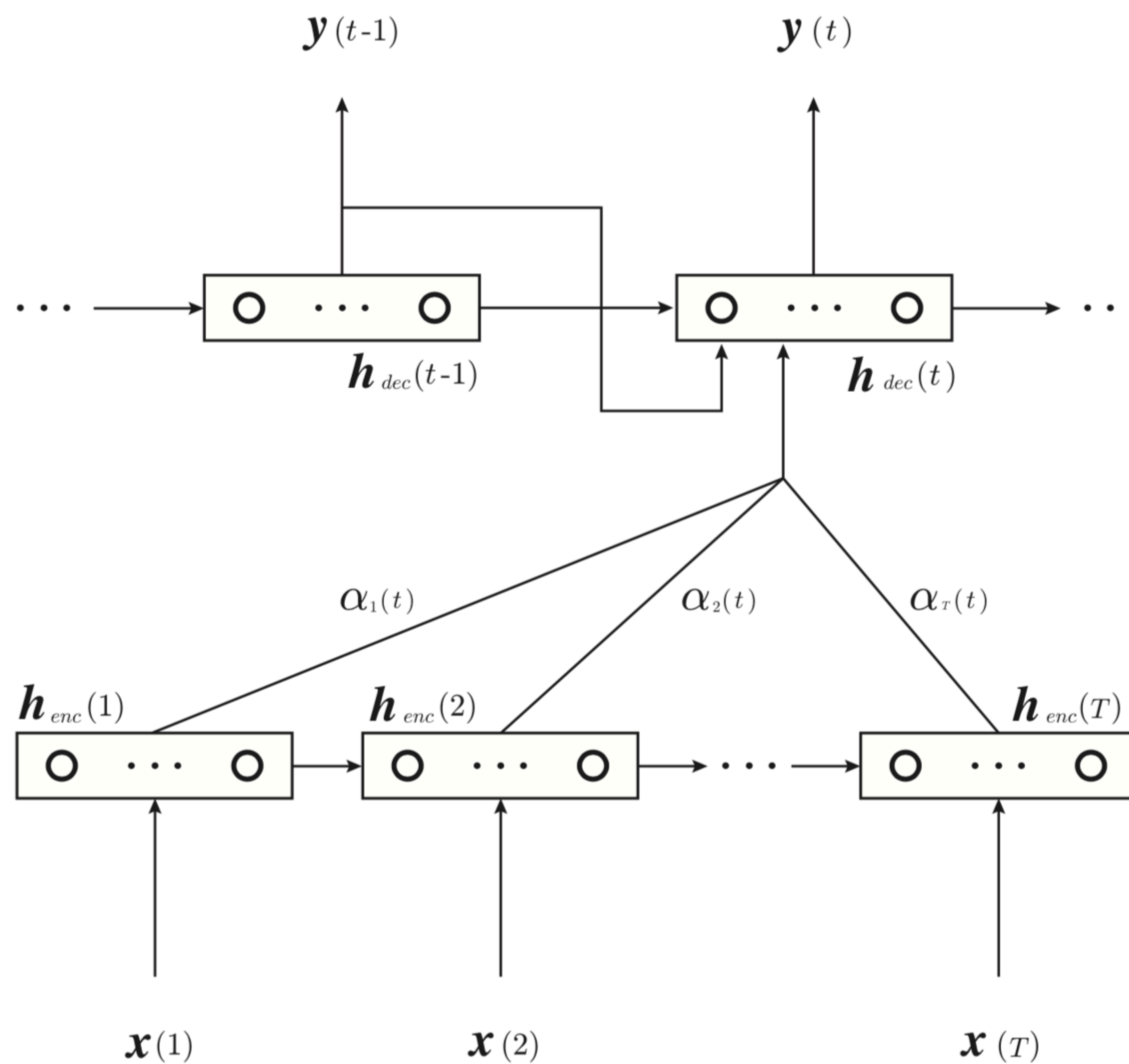
時間の重みを考慮し、各時刻によって動的に変わるベクトルを考えれば、
よりよいモデルが得られるはず

--> Attention

Encoder-Decoder



Attentionを取り入れた RNN Encoder-Decoder



$$\mathbf{h}_{dec}(t) = f(\mathbf{h}_{dec}(t-1), \mathbf{y}(t-1), \mathbf{c})$$



$$\mathbf{h}_{dec}(t) = f(\mathbf{h}_{dec}(t-1), \mathbf{y}(t-1), \underline{\mathbf{c}(t)})$$

時間を考慮

$$\mathbf{c}(t) = \sum_{\tau=1}^T \underline{\alpha_{\tau}(t)} \mathbf{h}_{enc}(\tau)$$

各エンコーダの値をどれだけデコーダに伝えるかの割合（＝ 重み）

$\alpha_\tau(t)$ は和が 1 になるので、

$w_\tau(t) := f(\mathbf{h}_{dec}(t-1), \mathbf{h}_{enc}(\tau))$ とおくと (=最適化すべき重み)



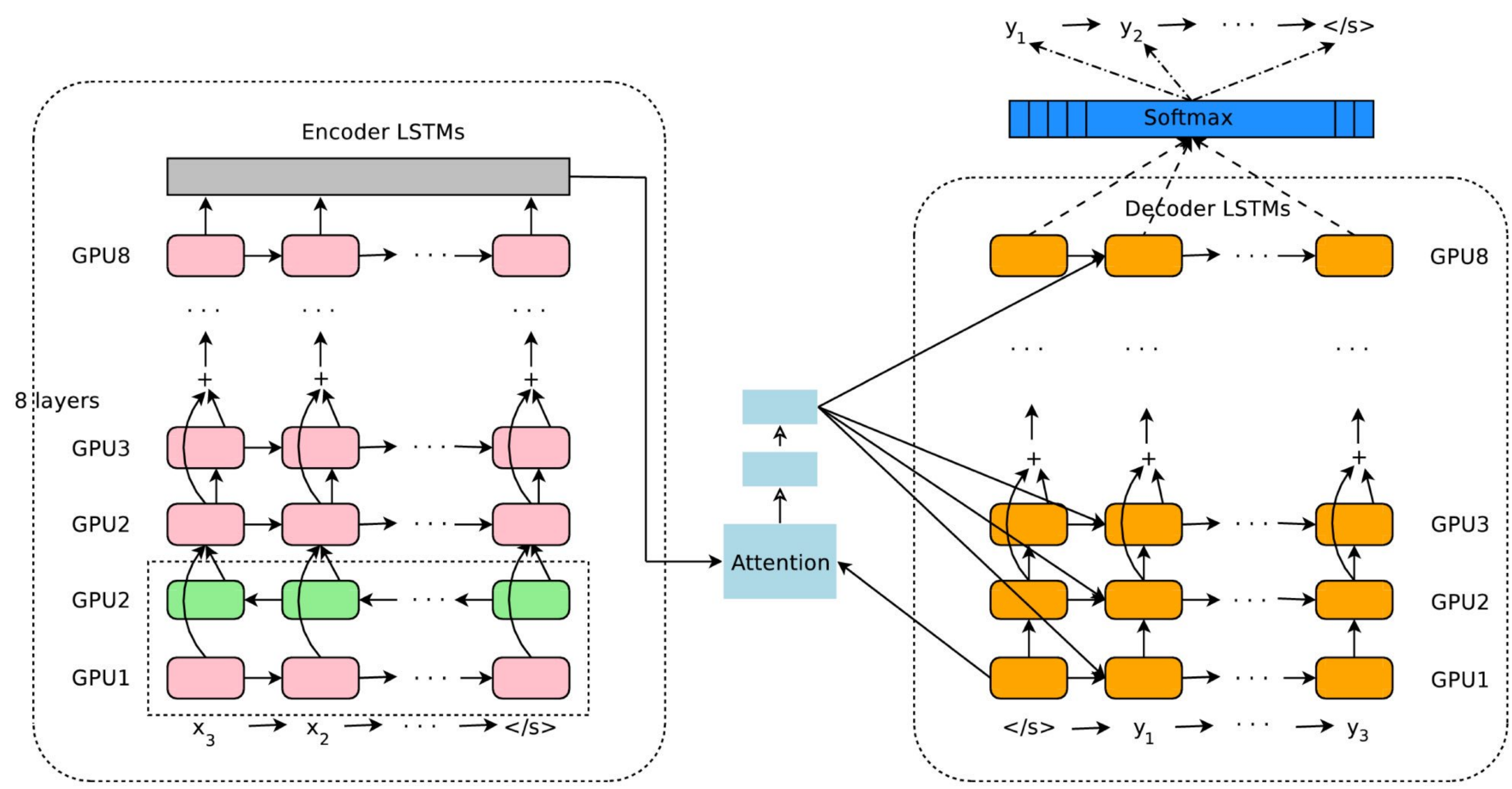
$$\alpha_\tau(t) = \frac{\exp(w_\tau(t))}{\sum_{\rho=1}^T \exp(w_\rho(t))} = \text{softmax}(w_\tau(t))$$

$$\begin{aligned}
 w_{\tau}(t) &= f(\mathbf{h}_{dec}(t-1), \mathbf{h}_{enc}(\tau)) \\
 &= \mathbf{v}_a^T \tanh(W_a \mathbf{h}_{dec}(t-1) + U_a \mathbf{h}_{enc}(\tau))
 \end{aligned}$$

ただし、重みに他の関数を用いるケースや、
Attentionの作り方にも別パターンがある

(再掲)

cf. Google翻訳



参考：LSTMにおけるAttention

LSTMブロックの出力

$$h(t) = f(h(t-1), x(t))$$

→ この式も、直前の時刻 $t-1$ にこれまでのすべての
時系列情報が集約されている形になっている

セル（CEC）自体に、時間の重みをもたせられないか？

→ 覗き穴結合なしのLSTMで考えてみる

$$\begin{pmatrix} \mathbf{i}(t) \\ \mathbf{o}(t) \\ \mathbf{f}(t) \\ \mathbf{a}(t) \end{pmatrix} = \begin{bmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{bmatrix} \cdot \begin{pmatrix} W_i & U_i & \mathbf{b}_i \\ W_o & U_o & \mathbf{b}_o \\ W_f & U_f & \mathbf{b}_f \\ W_a & U_a & \mathbf{b}_a \end{pmatrix} \begin{pmatrix} \mathbf{x}(t) \\ \underline{\mathbf{h}(t-1)} \\ 1 \end{pmatrix}$$

それぞれ、時間の重み
を考慮した項へ変更

$$\mathbf{c}(t) = \mathbf{i}(t) \odot \mathbf{a}(t) + \mathbf{f}(t) \odot \underline{\mathbf{c}(t-1)}$$



$$\begin{pmatrix} \tilde{\mathbf{h}}(t) \\ \tilde{\mathbf{c}}(t) \end{pmatrix} = \sum_{\tau=1}^{t-1} \alpha_{\tau}(t) \begin{pmatrix} \mathbf{h}(\tau) \\ \mathbf{c}(\tau) \end{pmatrix}$$

$$\alpha_{\tau}(t) = \text{softmax} (w_{\tau}(t))$$

$$w_{\tau}(t) \quad := \quad g \left(\boldsymbol{x}(t), \, \tilde{\boldsymbol{h}}(t - 1), \, \boldsymbol{h}(\tau) \right)$$

$$= \quad \boldsymbol{v}^T \tanh \left(W_x \boldsymbol{x}(t) + W_{\tilde{h}} \tilde{\boldsymbol{h}}(t - 1) + W_h \boldsymbol{h}(\tau) \right)$$

