$X_1, X_2, \ldots, X_n \longmapsto h$
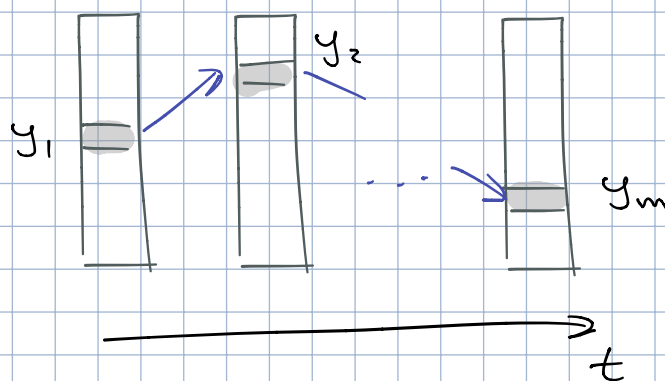
Generation process

$$P(y_1 | h)$$
$$P(y_2 | y_1, h)$$
$$\vdots$$

$$P(y_1 y_2 \cdots y_m | h) = P(y_1 | h) P(y_2 | y_1, h) \cdot \ldots \cdot P(y_m | y_1 \cdots y_{m-1} | h)$$

Each step of the way we predict prob dist over words



Therefore $P(y_{1 \cdots m} | h) = \hat{P}_{y_1} \cdots \hat{P}_{y_m}$

We want to find $\quad \underset{y_1 \cdots y_m}{\arg\max} \, P(y_1 \cdots y_m | h)$

At each timestamp $t$, we can keep fixed size pool of candidates, try to add new char, and leave only most likely sequences.
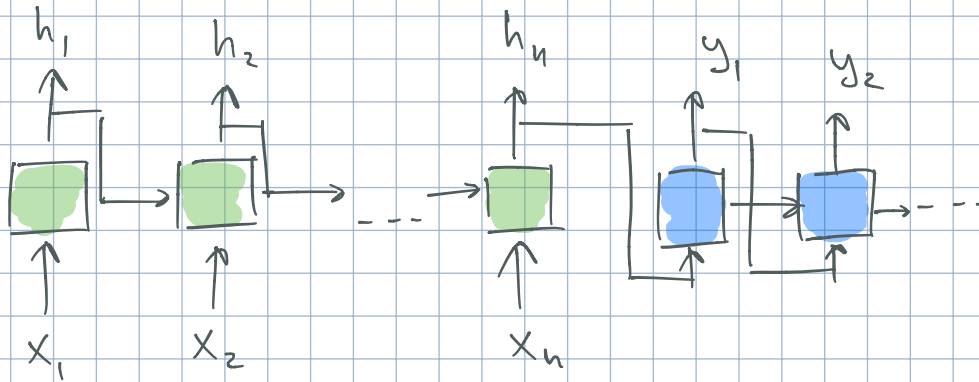
Fig: Encoder – Decoder Architecture

**Intuition** : Instead of using just last hidden state $(h = h_n)$ of encoder as an input to decoder at time t we can use

$$\alpha_1 h_1 + \alpha_2 h_2 + \cdots + \alpha_n h_n,$$

where $0 \le \alpha_i \le 1$ and $\sum \alpha_i = 1$,

For convinience $\alpha_i = \dfrac{\exp(\hat{\alpha_i})}{\sum\limits_{j} \exp(\hat{\alpha_j})}$, where

$\hat{\alpha_i}$ "depends" on position in the output $(y_t)$

More formaly, $\hat{\alpha_i} = MLP(h_i, h^d_{t-1})$, where $h^d_{t-1}$ — hidden state of decoder at previous timestamp.

# Transformer

## Self-Attention

$$A(q, K, V) = \text{attention based vcpr of a word}$$

**Intuition:**

Given words seq $x_1, x_2, \ldots, x_5$.

Vector representation for $x_3$ could be just simple emb. <u>OR</u> it could depend on the <u>context</u> $(x_1, \ldots, x_{i-1}, x_{i+1}, \ldots, x_n)$

$$A(q, k, V) = \sum_i' \frac{\exp(q \cdot k_i)}{\sum_j' \exp(q \cdot k_j)} \, v_i \quad,$$

where

$q$ — query $\quad (W^Q \cdot x_3)$

$k_i$ — key $\quad (W^k \cdot x_3)$ $\quad$ } Trainable matrices

$v_i$ — value $\quad (W^V \cdot x_3)$

we do this simulteneously for all words in a sentence and get $A_1, A_2, \ldots, A_5$

## Remark

$q_t$ — is like a question about $x_t$, and

$k_i$ — is an answer, there fore

$$\langle q_t, k_i \rangle - \text{measures how good}$$

$k_i$ answers $q_t$

Since we want weights to add up to 1
we use soft max

$$\alpha_i = \frac{\exp(q_t \cdot k_i)}{\sum_j \exp(q_t \cdot k_j)}$$

And to get final repr of $A_t$

$$\sum_i \alpha_i v_i$$

Matrix form of Self-Attention

$$A(Q, k, V) = \text{softmax}\left(\frac{Q k^T}{\sqrt{d_k}}\right) V$$

for scaling

row-wise softmax

$$Q = (q_1, q_2, \cdots, q_n)^T \in \mathbb{R}^{n \times d}$$
$$K = (k_1, k_2, \cdots, k_n)^T \in \mathbb{R}^{n \times d}$$
$$V = (v_1, v_2, \cdots, v_n)^T \in \mathbb{R}^{n \times d}$$
$$A(Q, K, V) = (A_1, A_2, \cdots, A_n)^T \in \mathbb{R}^{n \times d}$$

$A_i$ — context dependent vec repr of word $x_i$

## Multi-Head Attention

Recall that previously we were thinking about q's as questions

If we repeat comp of $A(Q, k, V)$ with
different matrices $W^Q, W^k, W^v$ it will give
us an opportunity to ask different questions

Let $h$ denote number of heads

$$head_i = A(W_i^Q \cdot Q, W_i^k \cdot k, W_i^v \cdot V), \quad i = 1, \dots, h$$

will give us $\mathbb{R}^{n \times d}$ self-attention matrices

Q: How to build final representation for each
word?

A:
$$\underbrace{concat(head_1, \dots, head_h)}_{\mathbb{R}^{n \times (d \cdot h)}} \cdot \underbrace{W_{out}}_{\mathbb{R}^{d \cdot h \times l}} \in \mathbb{R}^{n \times l}$$

✎ This approach allows us to build even
reacher context dependent words represent-s.

Transformer Architecture

Remarks
  * MLP after M-H Attention is applyed
independently for each position (aka row
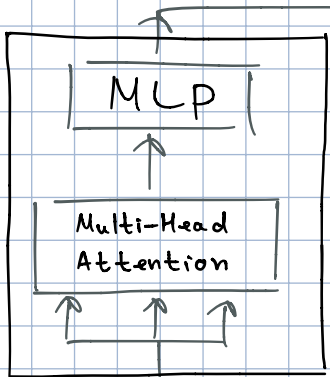of MHA matrix), hence outputs matrix

  * For stacking multiple layers, output
matrix from previous layer needs to be
transformed into matrices $Q, k, V$ (How?)

(Probably using simple mat mul)

Decoder

Softmax
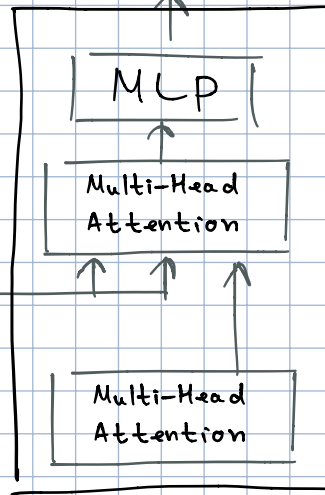
Linear

Encoder

MLP

Multi-Head
Attention

MLP

Multi-Head
Attention

Multi-Head
Attention

Stack
Nx
blocks

Stack
N
times

$Q, K, V$

$X_1, X_2, ---, X_n$

$y_1, y_2, ---$