

Foundations of Deep Learning



ALF

Alfredo Canziani

 @alfcnz

Attention (self/cross, hard/soft)

Dealing with sets

$$\mathbf{h} = \mathbf{X} \mathbf{a}$$

Self-attention (I)

$$\{\mathbf{x}_i\}_{i=1}^t = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_t\} \rightsquigarrow \mathbf{X} \in \mathbb{R}^{n \times t}, \quad \mathbf{x}_i \in \mathbb{R}^n$$

$$\mathbf{h} = \alpha_1 \mathbf{x}_1 + \alpha_2 \mathbf{x}_2 + \dots + \alpha_t \mathbf{x}_t = \mathbf{X} \mathbf{a} \in \mathbb{R}^n$$

$$\mathbf{X} \doteq \begin{bmatrix} | & | & & | \\ \mathbf{x}_1 & \mathbf{x}_2 & \dots & \mathbf{x}_t \\ | & | & & | \end{bmatrix} \in \mathbb{R}^{n \times t}$$

$$\begin{array}{l} \text{soft attention: } \rightarrow \|\mathbf{a}\|_1 = 1 \\ \text{hard attention: } \rightarrow \|\mathbf{a}\|_0 = 1 \end{array}$$

$$\textcolor{teal}{h} = \textcolor{violet}{X} \textcolor{brown}{a}$$

Self-attention (II)

$$\textcolor{brown}{a} = [\text{soft}](\arg\max_{\beta}(\textcolor{violet}{X}^{\top} \textcolor{violet}{x})) \in \mathbb{R}^t$$

$$\{\textcolor{violet}{x}_i\}_{i=1}^t \rightsquigarrow \{\textcolor{brown}{a}_i\}_{i=1}^t \rightsquigarrow \textcolor{brown}{A} \in \mathbb{R}^{t \times t}$$

$$\{\textcolor{brown}{a}_i\}_{i=1}^t \rightsquigarrow \{\textcolor{teal}{h}_i\}_{i=1}^t \rightsquigarrow \textcolor{teal}{H} \in \mathbb{R}^{n \times t}$$

$$\textcolor{teal}{H} = \textcolor{violet}{X} \textcolor{brown}{A} \in \mathbb{R}^{n \times t}$$

$\left[\cdot \right] : \text{optional}$

Key-value store

- Paradigm for
 - storing (saving)
 - retrieving (querying)
 - managing
- an associative array (dictionary / hash table)

Queries, keys, and values

$$\mathbf{q} = \mathbf{W}_{\mathbf{q}} \mathbf{x}, \quad \mathbf{k} = \mathbf{W}_{\mathbf{k}} \mathbf{x}, \quad \mathbf{v} = \mathbf{W}_{\mathbf{v}} \mathbf{x} \qquad \beta = \frac{1}{\sqrt{d}}$$

$$\mathbf{q}, \mathbf{k} \in \mathbb{R}^{d'}, \quad \mathbf{v} \in \mathbb{R}^{d''}, \quad d' = d'' \stackrel{\downarrow}{=} d$$

$$\{\mathbf{x}_i\}_{i=1}^t \rightsquigarrow \{\mathbf{q}_i\}_{i=1}^t, \{\mathbf{k}_i\}_{i=1}^t, \{\mathbf{v}_i\}_{i=1}^t \rightsquigarrow \mathbf{Q}, \mathbf{K}, \mathbf{V} \in \mathbb{R}^{d \times t}$$

$$\mathbf{a} = [\text{soft}](\arg\max_{\beta} (\mathbf{K}^{\top} \mathbf{q})) \in \mathbb{R}^t \qquad \mathbf{h} = \mathbf{V} \mathbf{a} \in \mathbb{R}^d$$

$$\{\mathbf{q}_i\}_{i=1}^t \rightsquigarrow \{\mathbf{a}_i\}_{i=1}^t \rightsquigarrow \mathbf{A} \in \mathbb{R}^{t \times t} \qquad \mathbf{H} = \mathbf{V} \mathbf{A} \in \mathbb{R}^{d \times t}$$

Implementation

$$\begin{bmatrix} q \\ k \\ v \end{bmatrix} = \begin{bmatrix} W_q \\ W_k \\ W_v \end{bmatrix} x \in \mathbb{R}^{3d}$$

$$\begin{bmatrix} q^1 \\ q^2 \\ \vdots \\ q^h \end{bmatrix} = \begin{bmatrix} W_q^1 \\ W_q^2 \\ \vdots \\ W_q^h \end{bmatrix} x = \begin{bmatrix} W_k^1 \\ W_k^2 \\ \vdots \\ W_k^h \end{bmatrix} x = \begin{bmatrix} v^1 \\ v^2 \\ \vdots \\ v^h \end{bmatrix} = \begin{bmatrix} W_v^1 \\ W_v^2 \\ \vdots \\ W_v^h \end{bmatrix} x$$

from the RNN lecture

$$h[t] = g(W_h [x^{[t]} \parallel h_{[t-1]}] + b_h)$$

$$h[0] \doteq \mathbf{0}, W_h \doteq \begin{bmatrix} W_{hx} & W_{hh} \end{bmatrix}$$

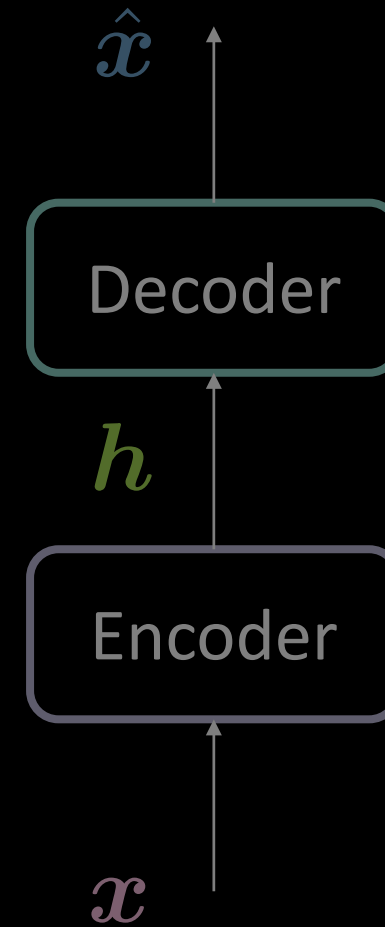
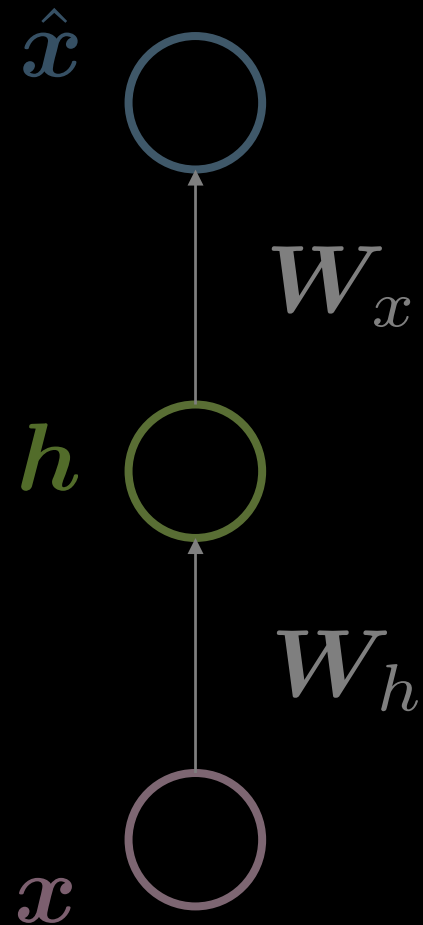
considering h heads we get a vector in \mathbb{R}^{3hd}

using a $W_h \in \mathbb{R}^{d \times hd}$ to go back to \mathbb{R}^d

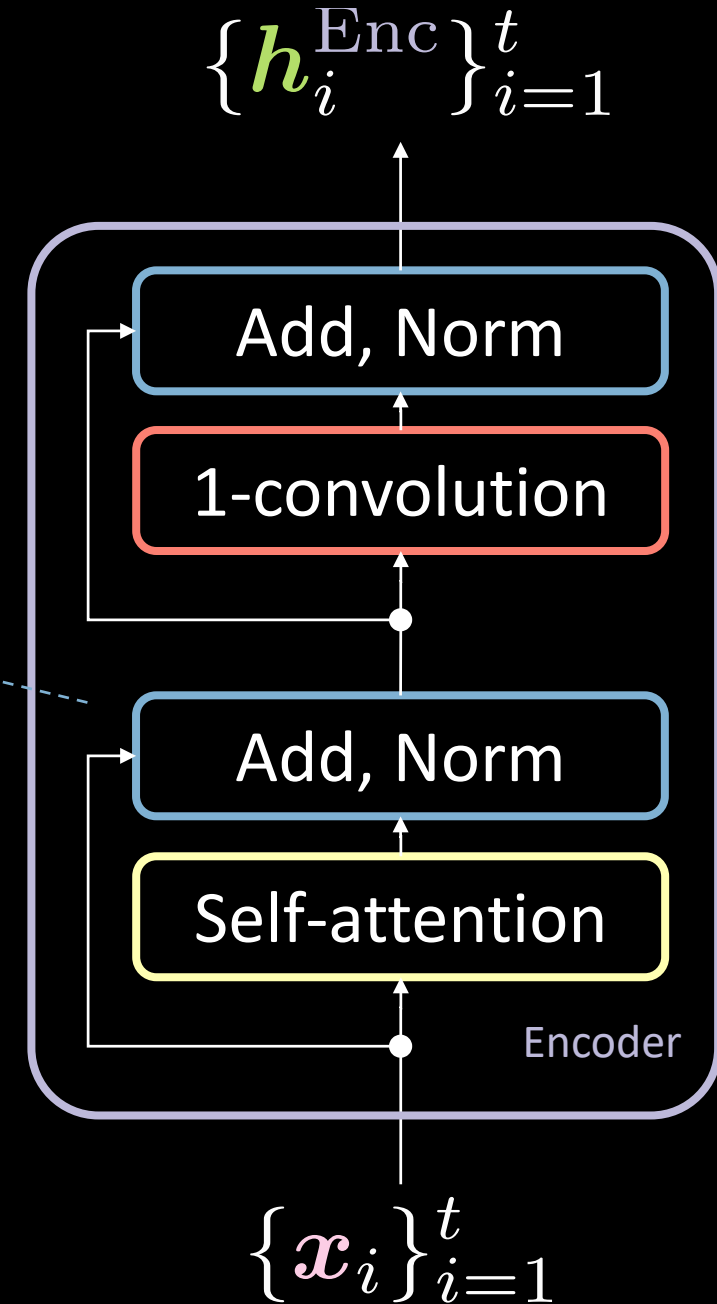
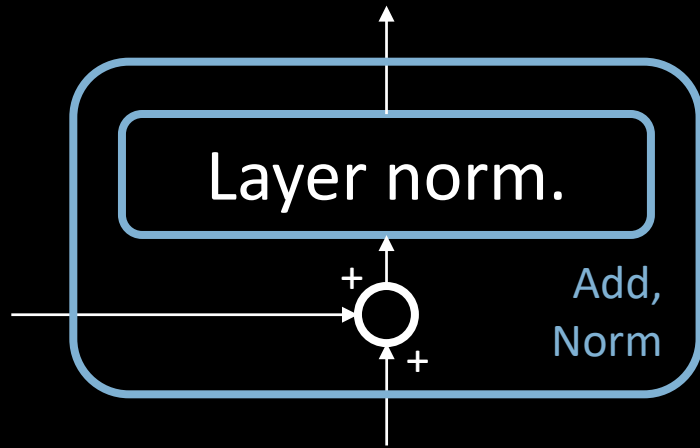
Transformer

Encoder-decoder architecture
(for Neural Machine Translation)

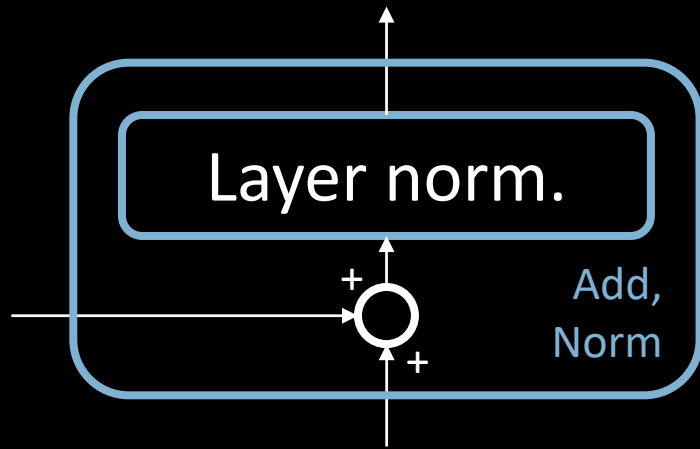
Auto-encoder (recap)



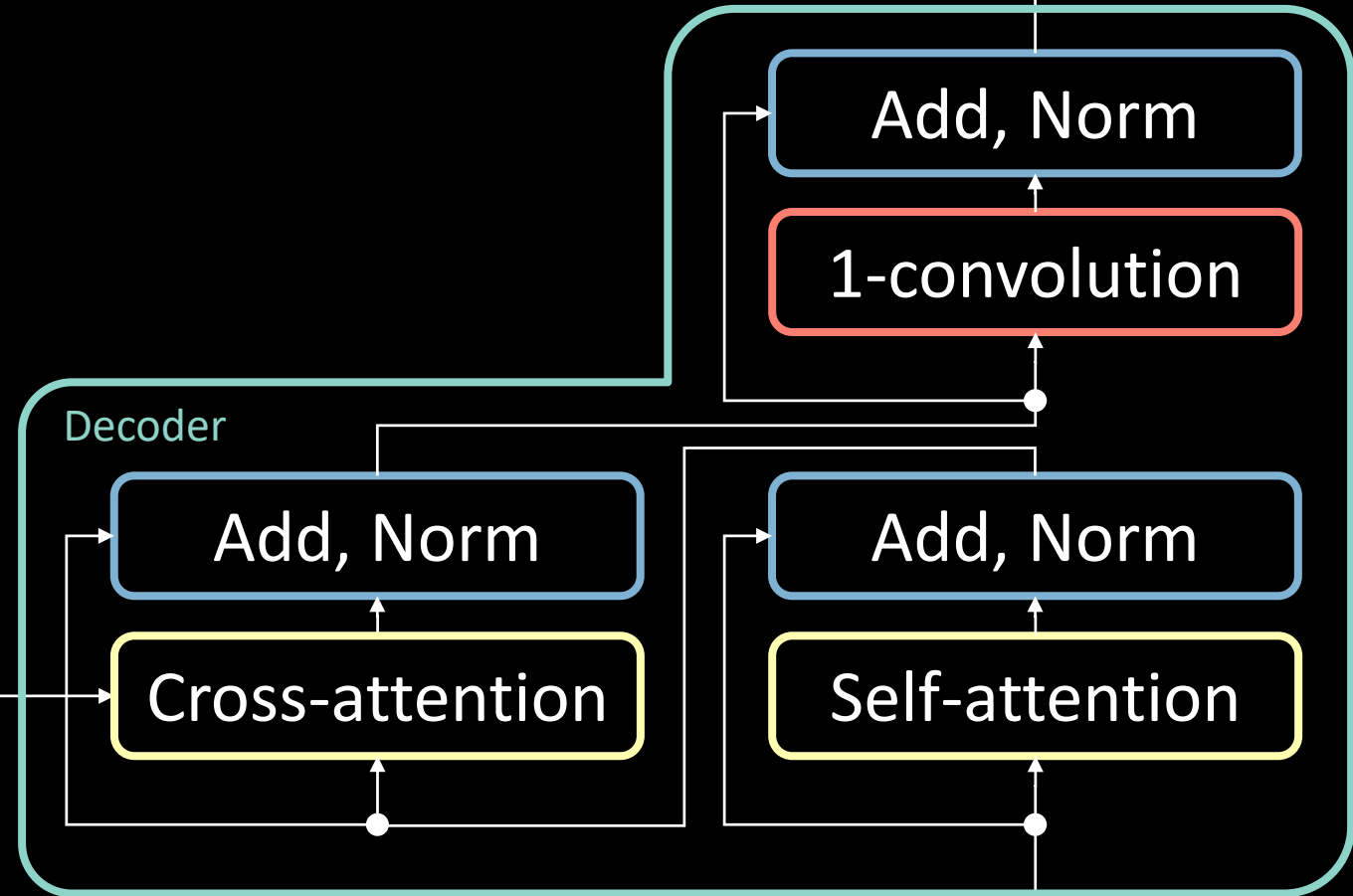
Transformer encoder



Transformer decoder



$\{h_i^{\text{Enc}}\}_{i=1}^t$



Decoder

$\{h_i^{\text{Dec}}\}_{i=1}^t$

Add, Norm

1-convolution

Add, Norm

Cross-attention

Add, Norm

Self-attention

$\{\hat{y}_i\}_{i=0}^{t-1}$