# Queries, keys, and values $\quad \{\boldsymbol{q}_i\}_{i=1}^{t} \rightsquigarrow \boldsymbol{Q} \in \mathbb{R}^{d' \times t}$

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

$$\boldsymbol{q}, \boldsymbol{k} \in \mathbb{R}^{d'}, \quad \boldsymbol{v} \in \mathbb{R}^{d''}$$

$$\{\boldsymbol{\xi}_j\}_{j=1}^{\tau} \rightsquigarrow \{\boldsymbol{k}_j\}_{j=1}^{\tau}, \{\boldsymbol{v}_j\}_{j=1}^{\tau} \rightsquigarrow \boldsymbol{K}, \boldsymbol{V} \in \mathbb{R}^{\{d', d''\} \times \tau}$$

$$\boldsymbol{a} = \text{softargmax}_{\beta}(\boldsymbol{K}^{\top} \boldsymbol{q}) \in \mathbb{R}^{\tau} \qquad \boldsymbol{h} = \boldsymbol{V} \boldsymbol{a} \in \mathbb{R}^{d''}$$

$$\{\boldsymbol{q}_i\}_{i=1}^{t} \rightsquigarrow \{\boldsymbol{a}_i\}_{i=1}^{t} \rightsquigarrow \boldsymbol{A} \in \mathbb{R}^{\tau \times t} \qquad \boldsymbol{H} = \boldsymbol{V} \boldsymbol{A} \in \mathbb{R}^{d'' \times t}$$

Self attention

$$d' = d'' \overset{\downarrow}{=} d$$

# Implementation

$$\boldsymbol{h}[t] = g\left(\boldsymbol{W_h}\begin{bmatrix}\boldsymbol{x}[t] \\ \boldsymbol{h}[t-1]\end{bmatrix} + \boldsymbol{b_h}\right)$$

$$\boldsymbol{h}[0] \doteq \boldsymbol{0}, \boldsymbol{W_h} \doteq [\boldsymbol{W_{hx}} \ \boldsymbol{W_{hh}}]$$

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

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

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

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

# Transformer

Encoders-predictor-decoder architecture
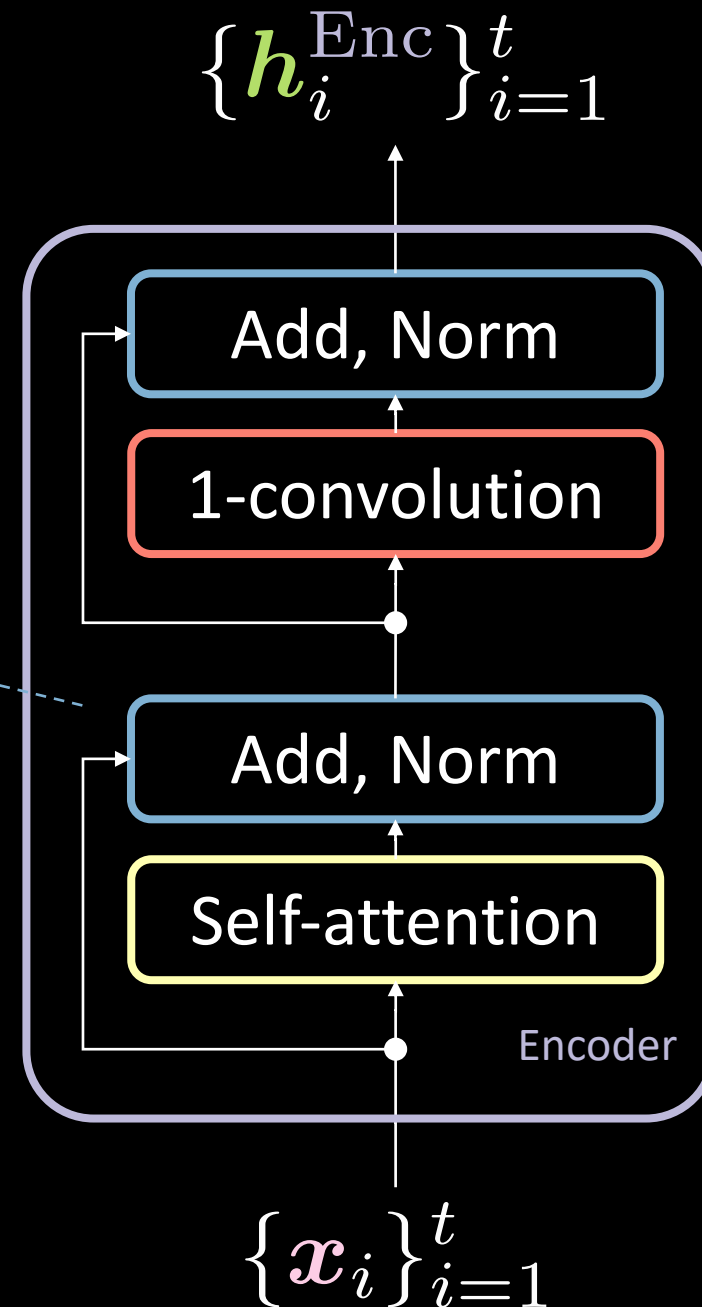
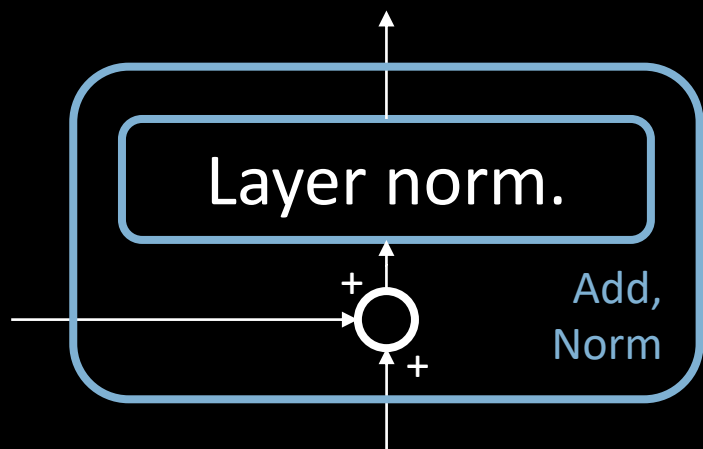(for Neural Machine Translation)

# Transformer encoder

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

Layer norm.

Add, Norm

Add, Norm

1-convolution

Add, Norm

Self-attention

Encoder

$$\{x_i\}_{i=1}^t$$

Transformer "decoder"

$\{\tilde{\boldsymbol{y}}_j\}_{j=1}^{\tau} \dashleftarrow \{\boldsymbol{h}_j^{\mathrm{Dec}}\}_{j=1}^{\tau}$

Layer norm.

Add, Norm

Add, Norm

1-convolution

"Decoder"

Add, Norm

Add, Norm

Cross-attention

Self-attention

$\{\boldsymbol{h}_i^{\mathrm{Enc}}\}_{i=1}^{t}$

inference $\{\tilde{\boldsymbol{y}}_j\}_{j=0}^{\tau-1}$     $\{\boldsymbol{y}_j\}_{j=0}^{\tau-1}$ training

# Transformer "decoder"

$$\{\tilde{\boldsymbol{y}}_j\}_{j=1}^{\tau} \dashleftarrow \{\boldsymbol{h}_j^{\mathrm{Dec}}\}_{j=1}^{\tau}$$

Layer norm.

Add, Norm

$$\{\boldsymbol{h}_i^{\mathrm{Enc}}\}_{i=1}^{t}$$

Add, Norm

1-convolution

Decoder

Add, Norm

Cross-attention

Predictor

Add, Norm

Self-attention

Encoder

inference $\{\tilde{\boldsymbol{y}}_j\}_{j=0}^{\tau-1}$ $\{\boldsymbol{y}_j\}_{j=0}^{\tau-1}$ training