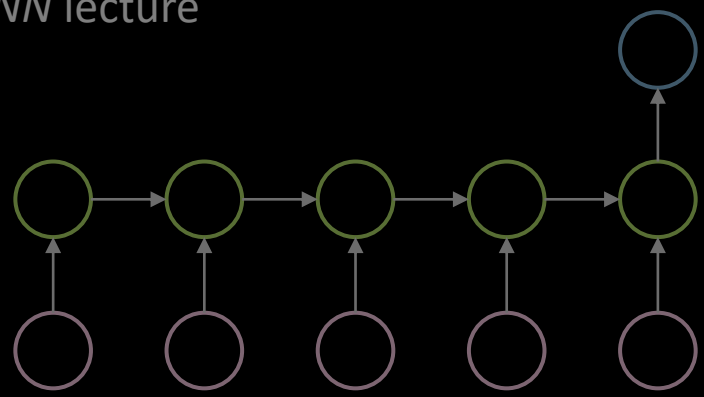


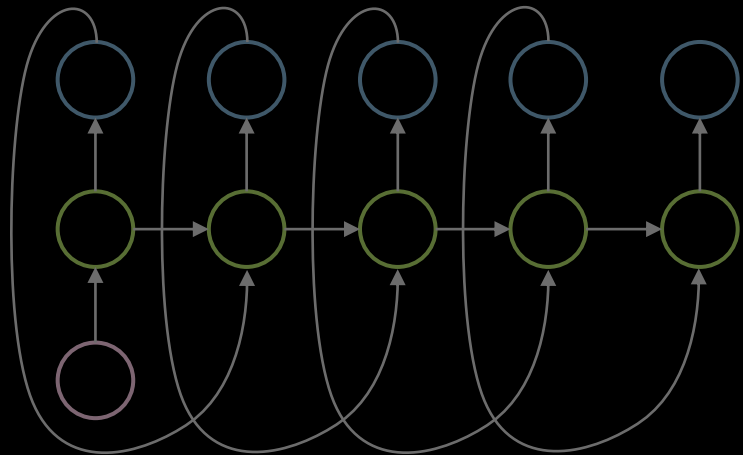
Attention (self/cross, hard/soft)

Dealing with sets

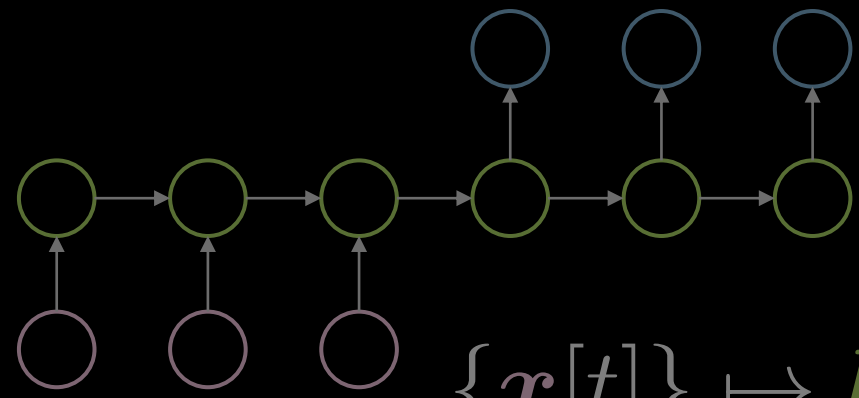
from the *RNN* lecture



$\{x[t]\} \mapsto \hat{y}[T] \quad \text{seq} \mapsto \text{vec}$

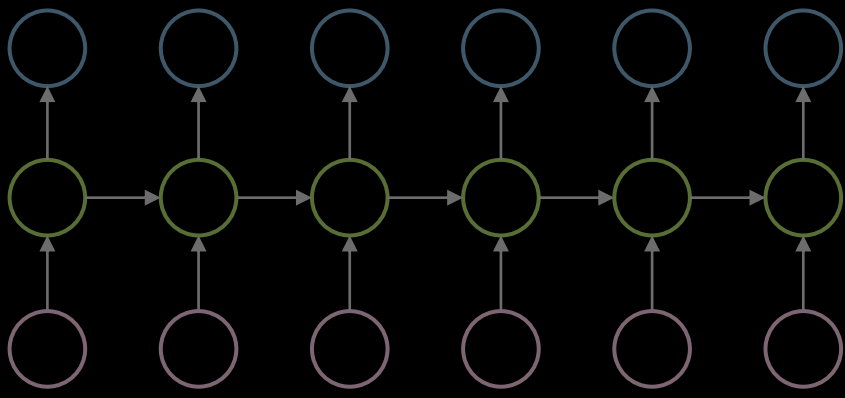


$x[1] \mapsto \{\hat{y}[t]\} \quad \text{vec} \mapsto \text{seq}$



$\{x[t]\} \mapsto h \mapsto \{\hat{y}[t]\}$

$\text{seq} \mapsto \text{vec} \mapsto \text{seq}$



$\{x[t]\} \mapsto \{\hat{y}[t]\} \quad \text{seq} \mapsto \text{seq}$

Use cases

- $\text{img} \mapsto \text{set}$: image to bounding box (BB) (DETER)
- $\text{set} \mapsto \text{set}$: point clouds to BB, surrounding vehicle trajectory pred.
- $\text{seq} \mapsto \text{seq}$: translation, conditional image generation (DALL-e)
- $\text{seq} \mapsto \text{set}$: event location and duration
- $\text{img} \mapsto \text{vec}$: visual image transformer (VIT)
- $\text{seq} \mapsto \text{vec}$: movies review

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

Self-attention (I)

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

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

$$\alpha_i > 0$$

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

$$\text{soft attention: } \left\{ \begin{array}{l} \|\mathbf{a}\|_1 = 1 \end{array} \right.$$

$$\text{hard attention: } \left\{ \begin{array}{l} \|\mathbf{a}\|_0 = 1 \end{array} \right.$$

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

Self-attention (II)

$$\mathbf{a} = \text{softmax}_{\beta}(\mathbf{X}^{\top} \mathbf{x}) \in \mathbb{R}^t$$

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

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

$$\mathbf{H} = \mathbf{X} \mathbf{A} \in \mathbb{R}^{n \times t}$$

Key-value store

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

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

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

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

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

$$\mathbf{a} = \text{softargmax}_{\beta}(\mathbf{K}^{\top} \mathbf{q}) \in \mathbb{R}^t \quad \mathbf{h} = \mathbf{V} \mathbf{a} \in \mathbb{R}^{d''}$$

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

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

$$\mathbf{q} = \mathbf{W}_q \mathbf{x}, \quad \mathbf{k} = \mathbf{W}_k \mathbf{x}, \quad \mathbf{v} = \mathbf{W}_v \mathbf{x} \quad \beta = \frac{1}{\sqrt{d'}}$$

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

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

$$\mathbf{a} = \text{softmax}_\beta(\mathbf{K}^\top \mathbf{q}) \in \mathbb{R}^\tau \quad \mathbf{h} = \mathbf{V} \mathbf{a} \in \mathbb{R}^{d''}$$

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

Self attention

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

Implementation

from the *RNN* lecture

$$\mathbf{h}[t] = g(\mathbf{W}_h [\mathbf{x}^{[t]} \mathbf{h}_{[t-1]}] + \mathbf{b}_h)$$

$$\mathbf{h}[0] \doteq \mathbf{0}, \mathbf{W}_h \doteq [\mathbf{W}_{hx} \mathbf{W}_{hh}]$$

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

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

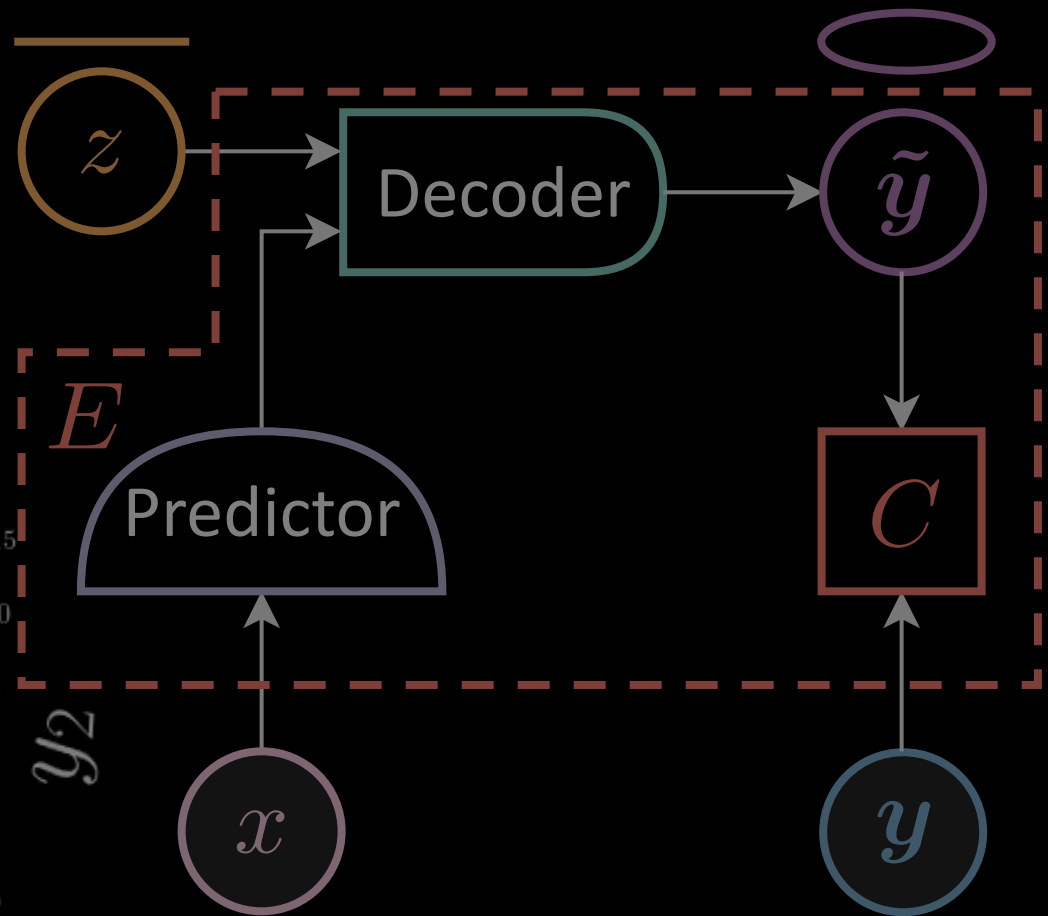
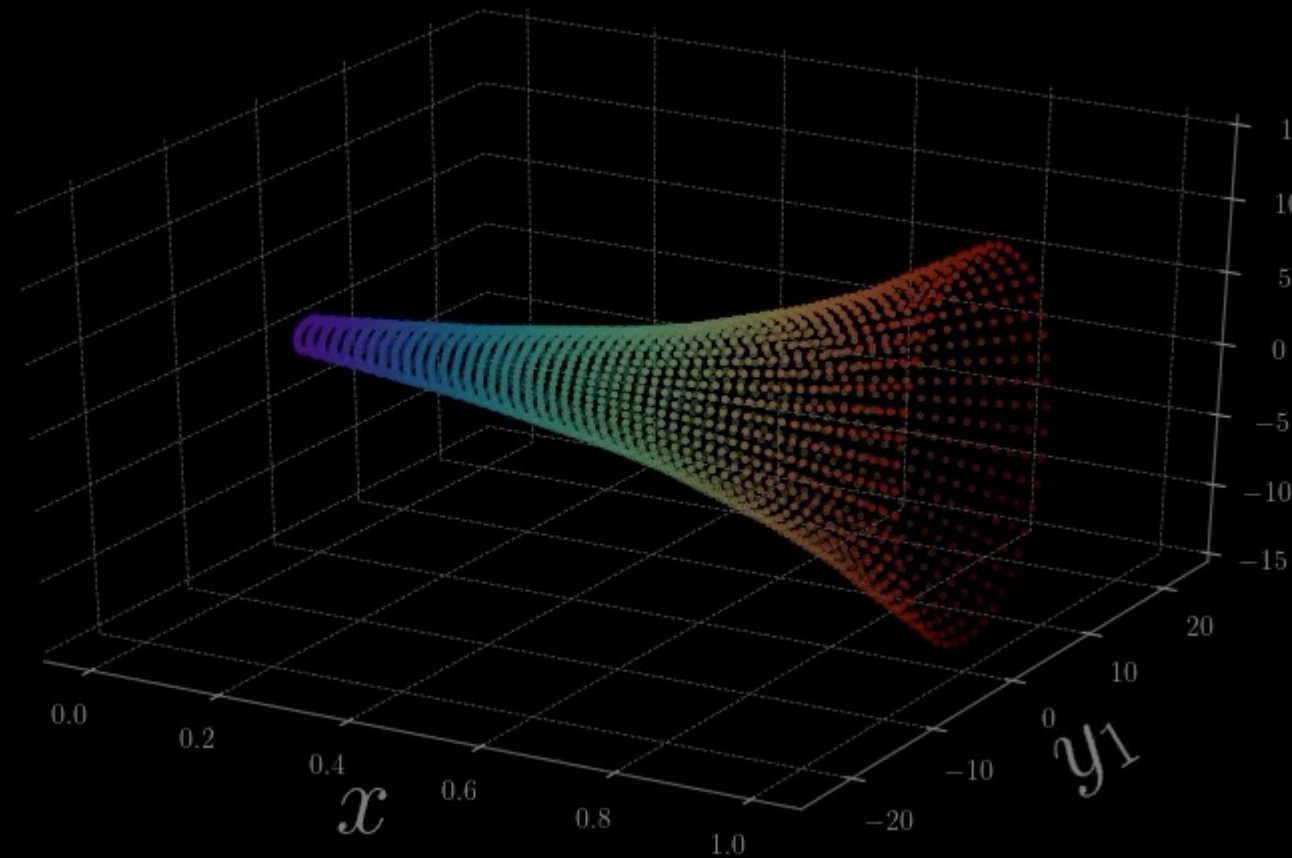
$$\begin{bmatrix} \mathbf{q}^1 \\ \mathbf{q}^2 \\ \vdots \\ \mathbf{q}^h \end{bmatrix} = \begin{bmatrix} \mathbf{W}_q^1 \\ \mathbf{W}_q^2 \\ \vdots \\ \mathbf{W}_q^h \end{bmatrix} \mathbf{x} = \begin{bmatrix} \mathbf{k}^1 \\ \mathbf{k}^2 \\ \vdots \\ \mathbf{k}^h \end{bmatrix} = \begin{bmatrix} \mathbf{W}_k^1 \\ \mathbf{W}_k^2 \\ \vdots \\ \mathbf{W}_k^h \end{bmatrix} \mathbf{x} = \begin{bmatrix} \mathbf{v}^1 \\ \mathbf{v}^2 \\ \vdots \\ \mathbf{v}^h \end{bmatrix} = \begin{bmatrix} \mathbf{W}_v^1 \\ \mathbf{W}_v^2 \\ \vdots \\ \mathbf{W}_v^h \end{bmatrix} \mathbf{x}$$

Transformer

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

from the *EBM* lecture

Trained model manifold



$$z = \left[0 : \frac{\pi}{24} : 2\pi \right]$$

$$x = \left[0 : \frac{1}{50} : 1 \right]$$

Autoencoder

$$\mathbf{h} = f(\mathbf{W}_h \mathbf{y} + \mathbf{b}_h)$$

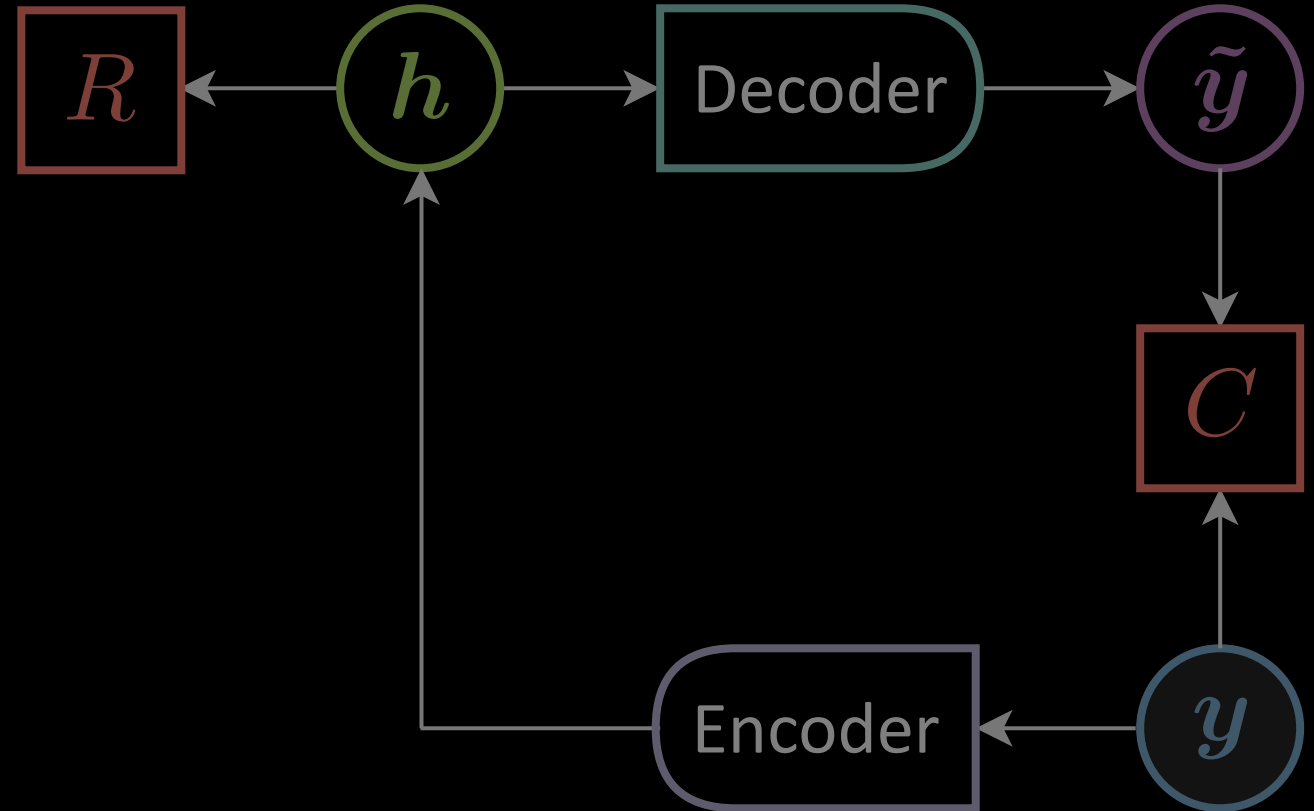
$$\tilde{\mathbf{y}} = g(\mathbf{W}_y \mathbf{h} + \mathbf{b}_y)$$

$$\mathbf{y}, \tilde{\mathbf{y}} \in \mathbb{R}^n$$

$$\mathbf{h} \in \mathbb{R}^d$$

$$\mathbf{W}_h \in \mathbb{R}^{d \times n}$$

$$\mathbf{W}_y \in \mathbb{R}^{n \times d}$$



$$\mathbf{h} = \text{Enc}(\mathbf{y}) \quad \tilde{\mathbf{y}} = \text{Dec}(\mathbf{h})$$

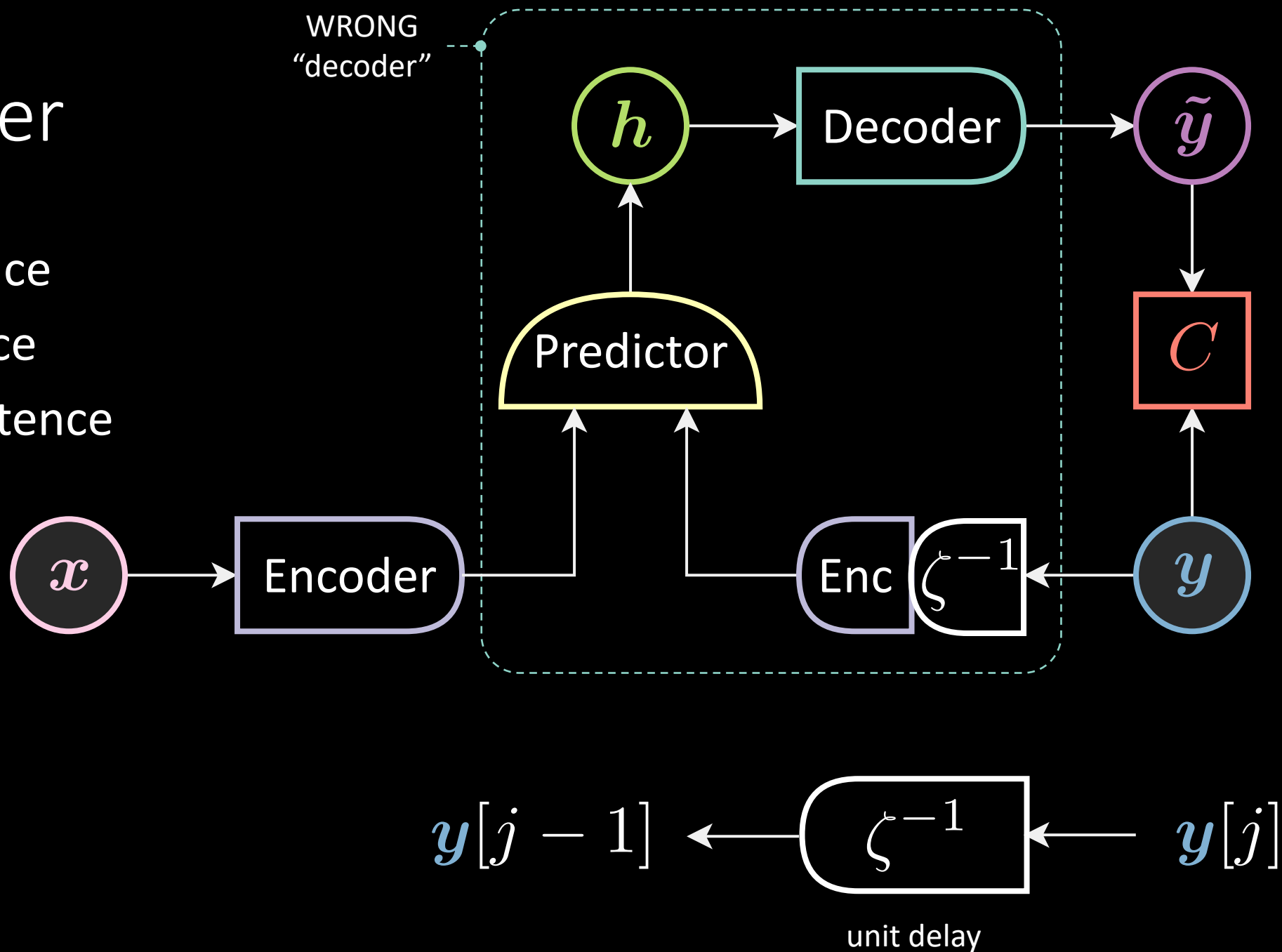
$$F(\mathbf{y}) = C(\mathbf{y}, \tilde{\mathbf{y}}) + R(\mathbf{h})$$

Transformer

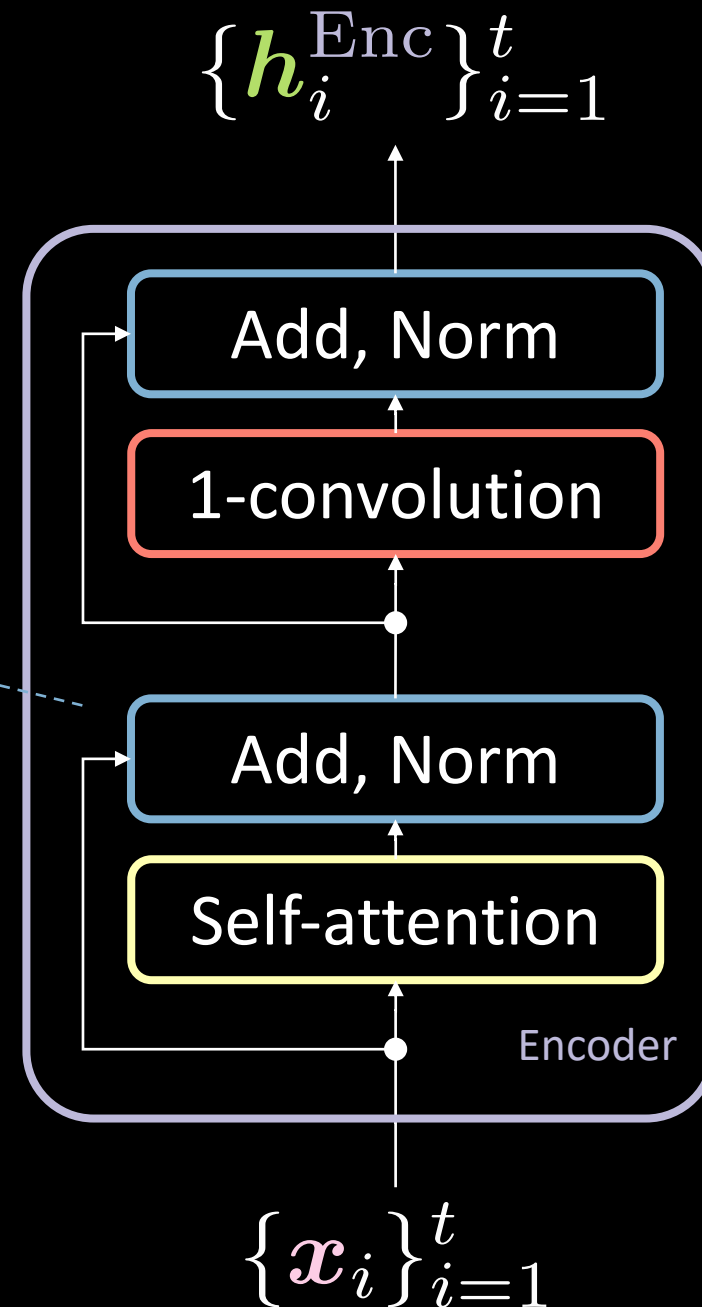
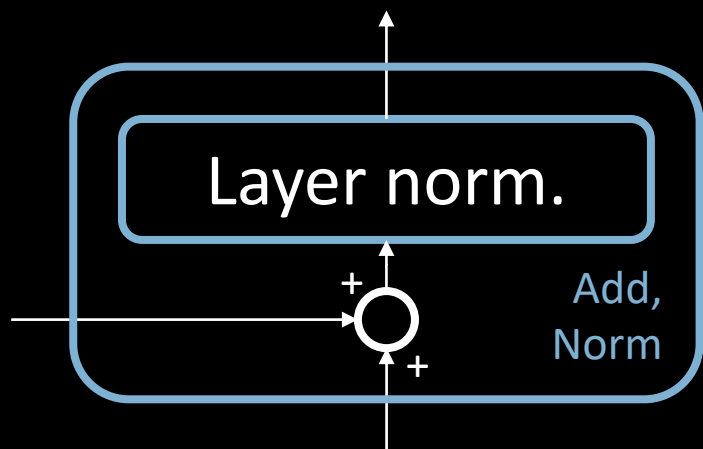
x source sentence

y target sentence

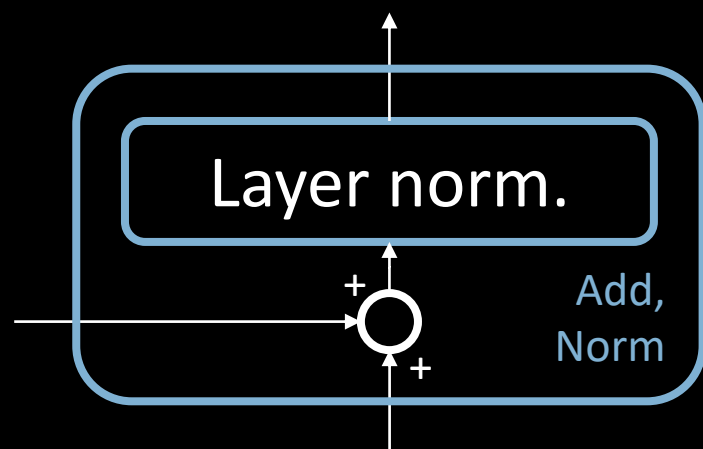
\tilde{y} predicted sentence



Transformer encoder



Transformer “decoder”



$\{\mathbf{h}_i^{\text{Enc}}\}_{i=1}^t$

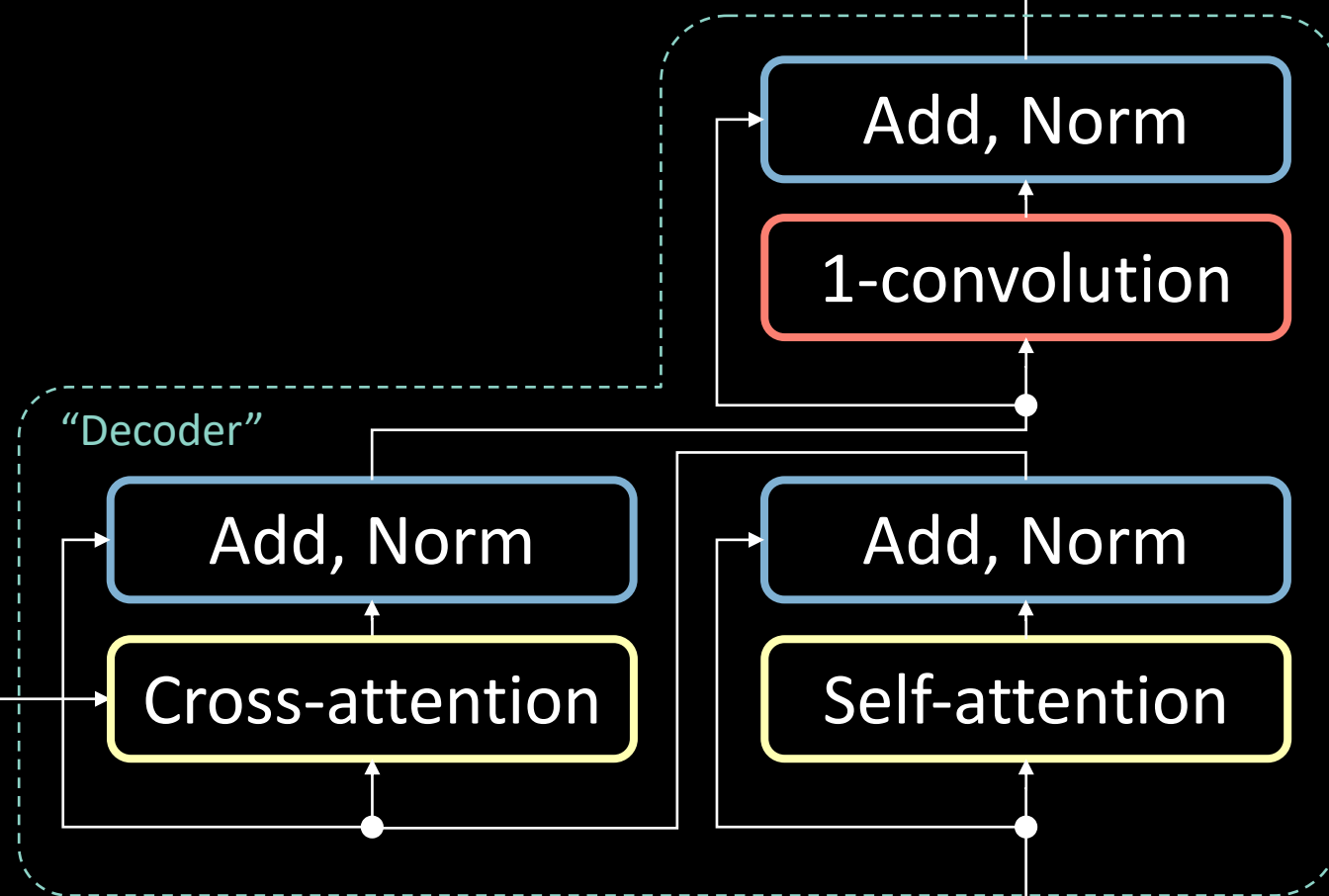
“Decoder”

inference

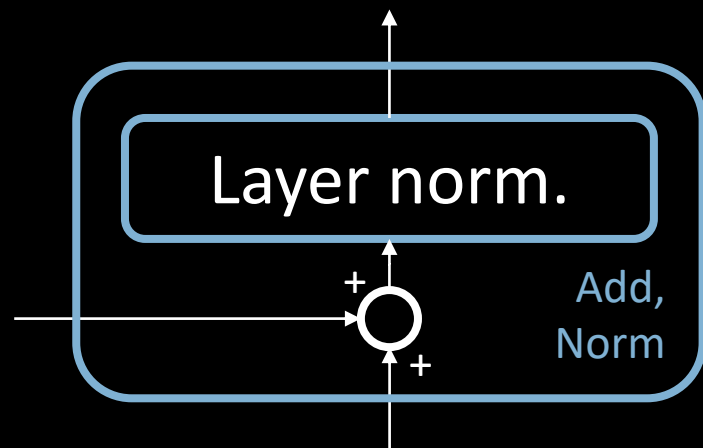
$\{\tilde{\mathbf{y}}_j\}_{j=0}^{\tau-1}$

$\{\mathbf{y}_j\}_{j=0}^{\tau-1}$

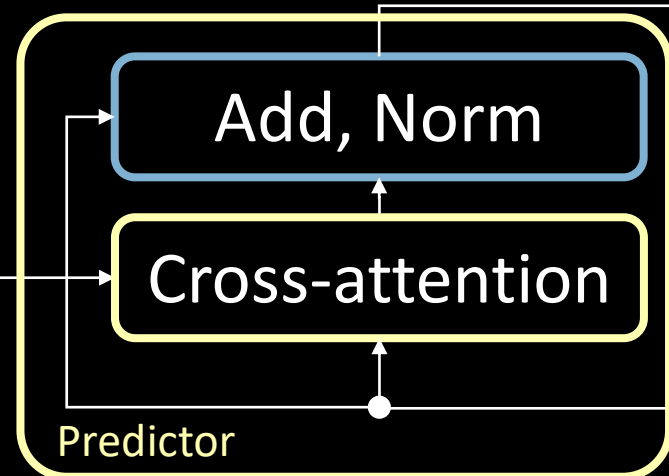
training



Transformer “decoder”

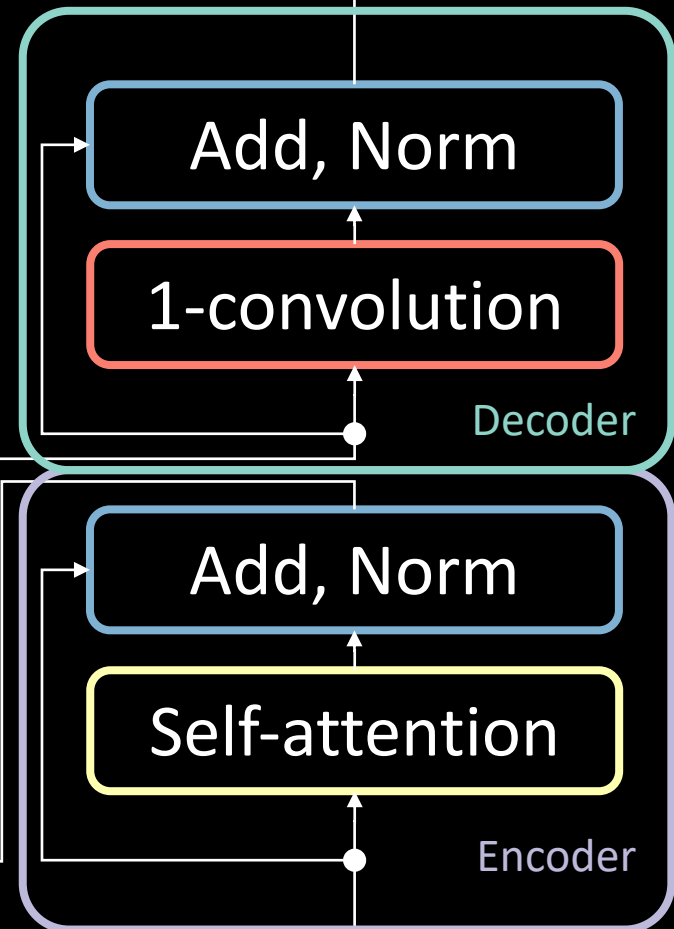


$\{h_i^{\text{Enc}}\}_{i=1}^t$



inference

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



training

$\{y_j\}_{j=0}^{\tau-1}$

$\{\tilde{y}_j\}_{j=1}^{\tau} \leftarrow \{h_j^{\text{Dec}}\}_{j=1}^{\tau}$