

ECE/SIOC 228 Machine Learning for Physical Applications

Lecture 7: Recurrent Neural Networks

Yuanyuan Shi

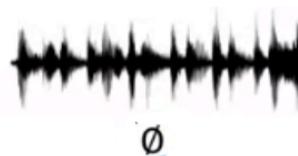
Assistant Professor, ECE

University of California, San Diego

Examples of Sequence Data

UC San Diego

Speech recognition



→ “The quick brown fox jumped over the lazy dog.”

Music generation



Sentiment classification

“There is nothing to like
in this movie.”



DNA sequence analysis

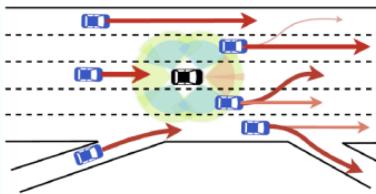
AGCCCCCTGTGAGGAACCTAG

AG~~CCCC~~TGTGAGGAACCTAG

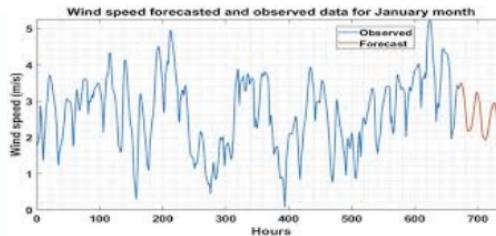
Machine translation

Voulez-vous chanter avec
moi?

Do you want to sing with
me?



trajectory prediction



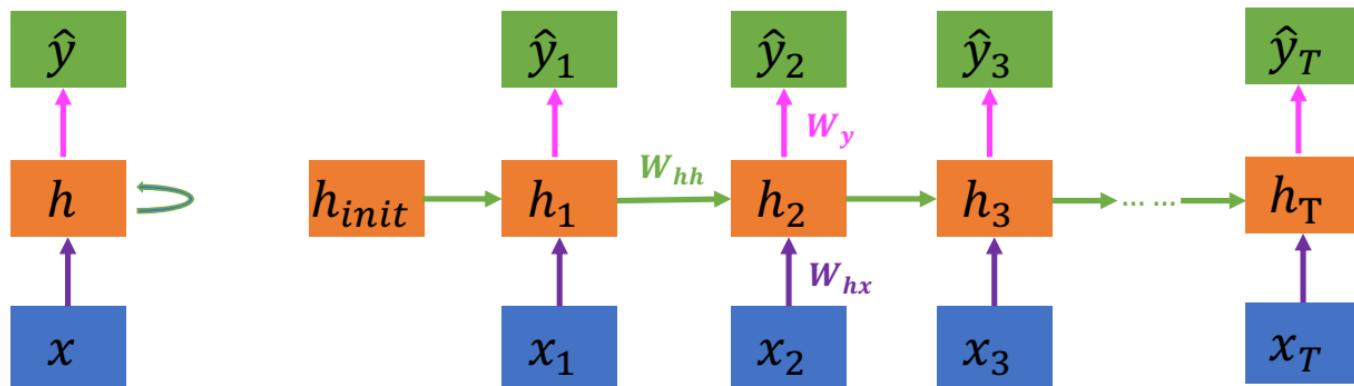
Outline

UC San Diego

- Recurrent Neural Networks
- Modern Recurrent Neural Networks
- Attention Mechanism and Transformer

Recurrent Neural Networks

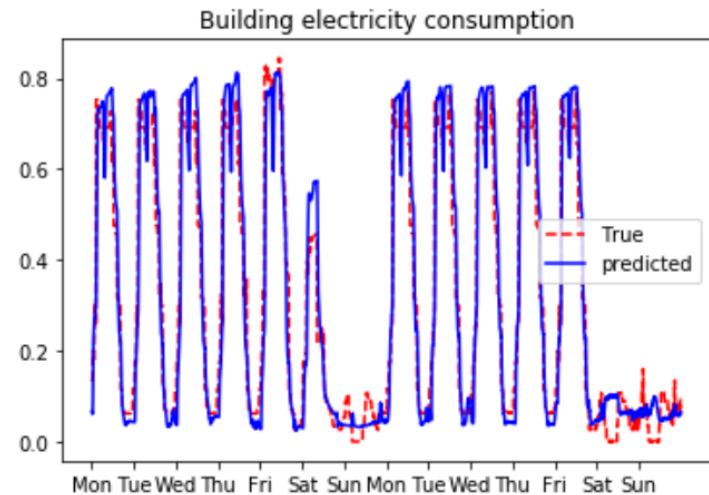
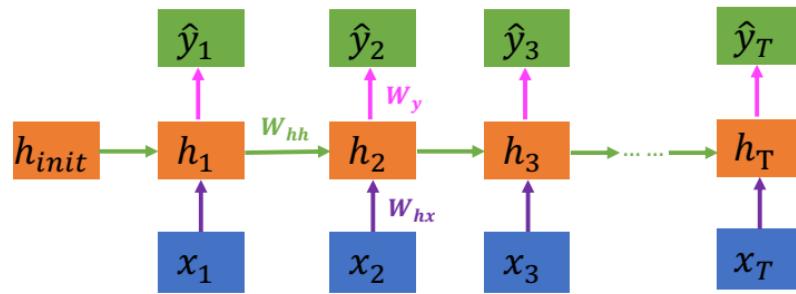
UC San Diego



- Computation takes into historical info $\rightarrow t$
- Weights are shared across time each set of colored arrows denotes the same weights
- Process input of any length ; model size not increase with size of input

Recurrent Neural Networks: Example

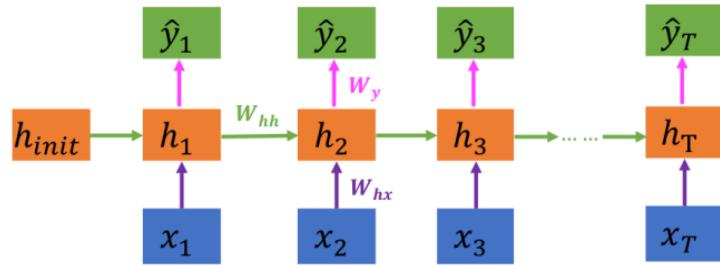
UC San Diego



- \hat{y}_t denotes the predicted electricity consumption, at time t
- x_t denotes the features that affects the building energy consumption, e.g., environment features (temperature, humidity, ...), occupancy, room temperature setpoint...
- h_t encodes the effect of historical $[x_1, \dots, x_{t-1}]$

Forward Computation

UC San Diego



for $t=1, \dots, T$

$$h_t = \sigma(W_{hh} h_{t-1} + W_{hx} x_t + b_h)$$

$$\hat{y}_t = \sigma(W_y h_t + b_y)$$

Loss function

$$L = \frac{1}{T} \sum_{t=1}^T L(\hat{y}_t, y_t)$$

Suppose we have n training samples:

$$x_t \in R^{d \times n}, h_t \in R^{h \times n}, \hat{y}_t \in R^{q \times n}$$

$$W_{hx} \in R^{h \times d}$$

$$W_{hh} \in R^{h \times h}$$

$$W_y \in R^{q \times h}$$

Backpropagation

UC San Diego

Two simplification: 1) ignore bias b_h, b_y ; 2) set activation $\delta(z) = z$

$$\frac{\partial L}{\partial \hat{y}_t} = \frac{1}{T} \frac{\partial L(y_t, \hat{y}_t)}{\partial \hat{y}_t} \in R^{q \times n}$$

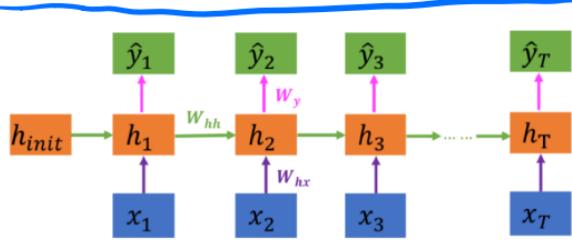
Now we compute gradient w.r.t. W_y (parameter
of the output layer)

L depends on W_y via $\hat{y}_1, \hat{y}_2, \dots, \hat{y}_T$

$$\frac{\partial L}{\partial W_y} = \sum_{t=1}^T \frac{\partial L}{\partial \hat{y}_t} \cdot \frac{\partial \hat{y}_t}{\partial W_y} = \sum_{t=1}^T \frac{\partial L}{\partial \hat{y}_t} h_t^T$$

Backpropagation

UC San Diego



for $t=1, \dots, T$

$$h_t = \sigma(W_{hh} h_{t-1} + W_{hx} x_t + b_h)$$

$$\hat{y}_t = \sigma(W_y h_t + b_y)$$

Loss function

$$L = \frac{1}{T} \sum_{t=1}^T L(\hat{y}_t, y_t)$$

Now, we look at gradient w.r.t. h_t

$$\frac{\partial L}{\partial h_T} = \text{prod}\left(\frac{\partial L}{\partial \hat{y}_T}, \frac{\partial \hat{y}_T}{\partial h_T}\right) = W_y^T \frac{\partial L}{\partial \hat{y}_T}$$

For $t < T$, we want to compute $\frac{\partial L}{\partial h_t}$ suppose $\frac{\partial L}{\partial h_{t+1}}$ and $\frac{\partial L}{\partial \hat{y}_t}$ are computed

$$\begin{aligned} \frac{\partial L}{\partial h_t} &= \text{prod}\left(\frac{\partial L}{\partial h_{t+1}}, \frac{\partial h_{t+1}}{\partial h_t}\right) + \text{prod}\left(\frac{\partial L}{\partial \hat{y}_t}, \frac{\partial \hat{y}_t}{\partial h_t}\right) \\ &= W_{hh}^T \frac{\partial L}{\partial h_{t+1}} + W_y^T \frac{\partial L}{\partial \hat{y}_t} \end{aligned}$$

Backpropagation

UC San Diego

$$\begin{aligned}\frac{\partial L}{\partial h_t} &= W_{hh}^T \frac{\partial L}{\partial h_{t+1}} + W_y^T \frac{\partial L}{\partial \hat{y}_t} \\&= W_{hh}^T \left(W_{hh}^T \frac{\partial L}{\partial h_{t+2}} + W_y^T \frac{\partial L}{\partial \hat{y}_{t+1}} \right) + W_y^T \frac{\partial L}{\partial \hat{y}_t} \\&= (W_{hh}^T)^2 \cdot \frac{\partial L}{\partial h_{t+2}} + W_{hh}^T W_y^T \frac{\partial L}{\partial \hat{y}_{t+1}} + W_y^T \frac{\partial L}{\partial \hat{y}_t} \\&= (W_{hh}^T)^2 \cdot \left[W_{hh}^T \frac{\partial L}{\partial \hat{y}_{t+3}} + W_y^T \frac{\partial L}{\partial \hat{y}_{t+2}} \right] + W_{hh}^T W_y^T \frac{\partial L}{\partial \hat{y}_{t+1}} + W_y^T \frac{\partial L}{\partial \hat{y}_t} \\&= \dots = \sum_{i=t}^T (W_{hh}^T)^{T-i} W_y^T \frac{\partial L}{\partial \hat{y}_{T+t-i}}\end{aligned}$$

Backpropagation

UC San Diego

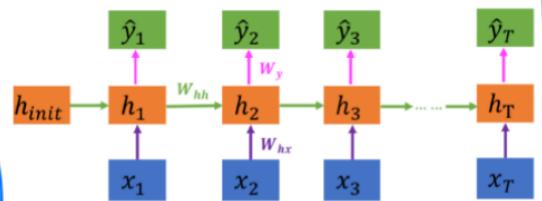
$$\frac{\partial L}{\partial h_t} = \sum_{i=t}^T (W_{hh}^T)^{T-i} w_y^T \frac{\partial L}{\partial \hat{y}_{T+t-i}}$$

- It involves (potentially) very large power of W_{hh}^T
- If $\text{eigen}(W_{hh}^T) > 1$, diverge (gradient exploding)
- If $\text{eigen}(W_{hh}^T) < 1$, very small (gradient vanishing)

★ Truncate the time steps $\frac{\partial L}{\partial h_t} = \sum_{i=t}^{t+H} (W_{hh}^T)^{t+H-i} w_y^T \frac{\partial L}{\partial \hat{y}_{t+H+i}}$

Backpropagation

UC San Diego



for $t=1, \dots, T$

$$h_t = \sigma(W_{hh} h_{t-1} + W_{hx} x_t + b_h)$$
$$\hat{y}_t = \sigma(W_y h_t + b_y)$$

Loss function

$$L = \frac{1}{T} \sum_{t=1}^T L(\hat{y}_t, y_t)$$

$$\frac{\partial L}{\partial W_{hx}} = \sum_{t=1}^T \text{prod} \left(\frac{\partial L}{\partial h_t}, \frac{\partial h_t}{\partial W_{hx}} \right) = \sum_{t=1}^T \frac{\partial L}{\partial h_t} x_t^T$$

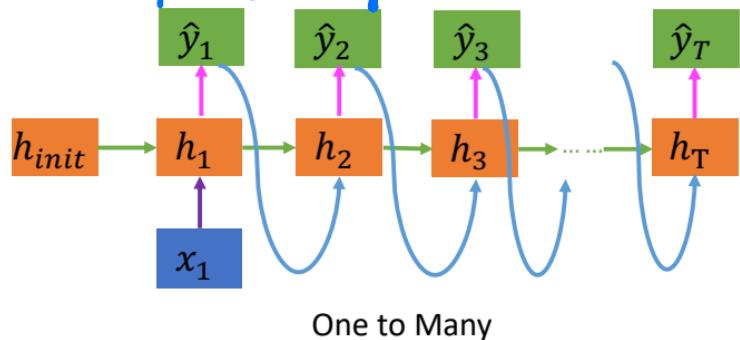
$$\frac{\partial L}{\partial W_{hh}} = \sum_{t=1}^T \text{prod} \left(\frac{\partial L}{\partial h_t}, \frac{\partial h_t}{\partial W_{hh}} \right) = \sum_{t=1}^T \frac{\partial L}{\partial h_t} h_{t-1}^T$$

- $\frac{\partial L}{\partial h_t}$ also affect the gradient magnitude $\frac{\partial L}{\partial W_{hx}}, \frac{\partial L}{\partial W_{hh}}$
- Training RNN alternates forward & backward propagation through time

Different Types of RNNs

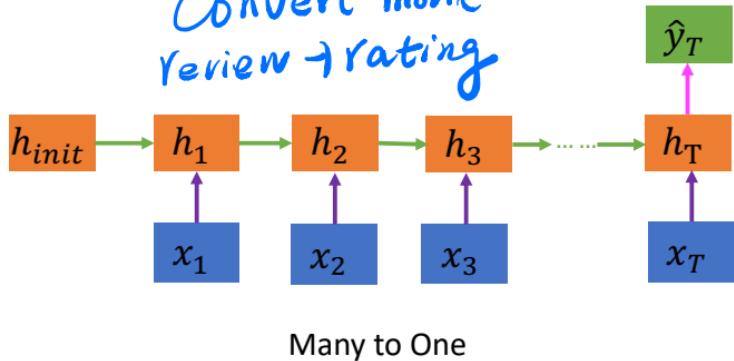
UC San Diego

MUSIC generation

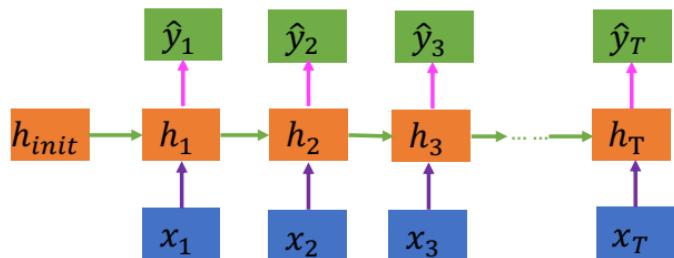


One to Many

Convert movie
review \rightarrow rating



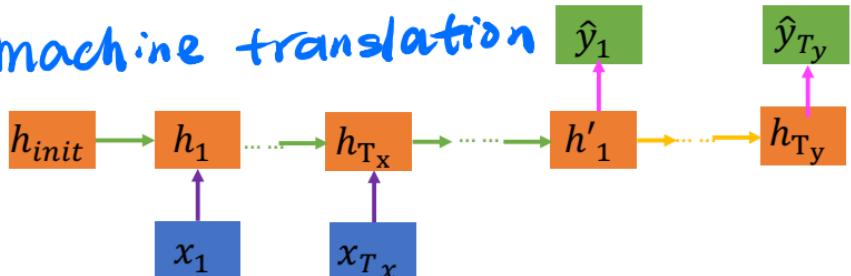
Many to One



Many to Many

building energy
prediction

machine translation



Many to Many (input and output have
different length)

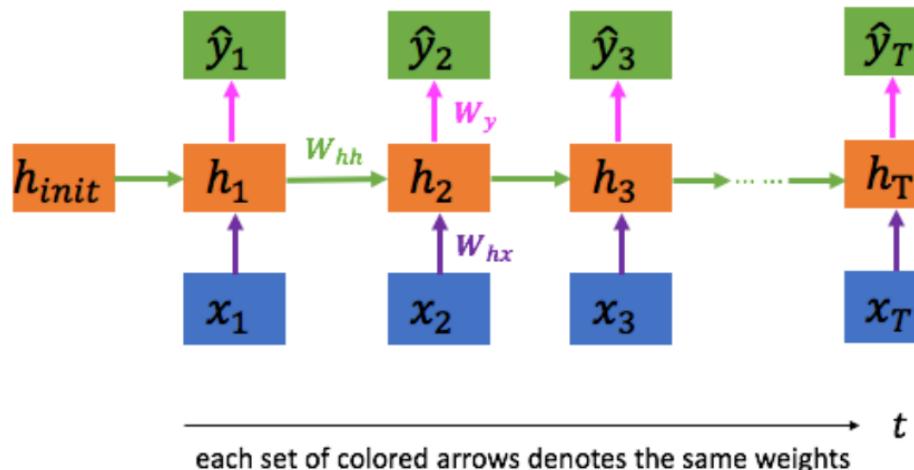
Outline

UC San Diego

- Recurrent Neural Networks
- **Modern Recurrent Neural Networks**
- Attention Mechanism and Transformer

Vanishing Gradient with RNN

UC San Diego

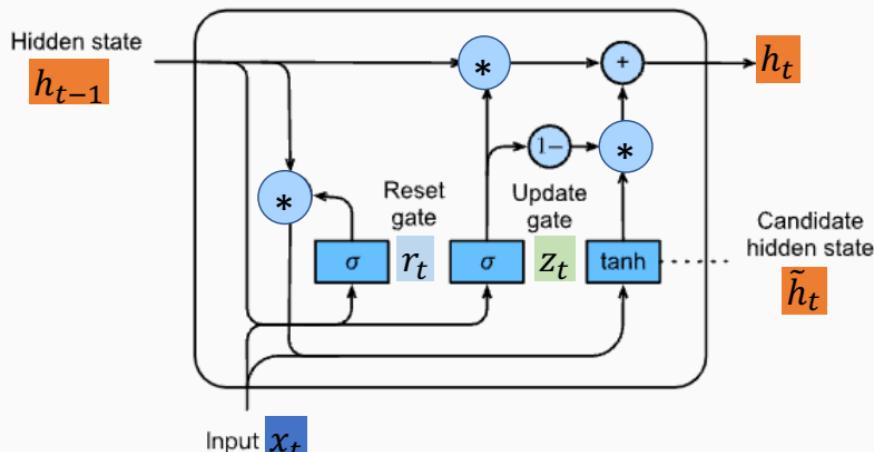


- Long products of matrices can lead to vanishing or exploding gradients

$$\frac{\partial L}{\partial h_t} = \sum_{i=t}^T (W_{hh}^T)^{T-i} W_y^T \frac{\partial L}{\partial \hat{y}_{T+t-i}}$$

Gated Recurrent Units (GRU)

UC San Diego



σ

FC layer with
activation function

Elementwise
operator

Copy

Concatenate

Reset gate:

$$r_t = \text{sigmoid}(W_{rx}x_t + W_{rh}h_{t-1} + b_r)$$

Update gate,

$$z_t = \text{sigmoid}(W_{zx}x_t + W_{zh}h_{t-1} + b_z)$$

$$x_t \in \mathbb{R}^{d \times n}, h_{t-1} \in \mathbb{R}^{h \times n},$$

$$r_t, z_t \in \underline{\mathbb{R}^{h \times n}}$$

$$W_{rx}, W_{zx} \in \mathbb{R}^{h \times d}$$

$$W_{rh}, W_{zh} \in \mathbb{R}^{h \times h}$$

Gated Recurrent Units (GRU)

UC San Diego

Reset gate: $R_t = \text{sigmoid}(W_{rx}x_t + W_{rh}h_{t-1} + b_r)$

Update gate: $Z_t = \text{sigmoid}(W_{zx}x_t + W_{zh}h_{t-1} + b_z)$

proposal / candidate hidden state:

$$\tilde{h}_t = \tanh(W_h x_t + W_{hh} \underset{\uparrow}{(R_t * h_{t-1})} + b_h)$$

element-wise

- If $R_t \approx 1$, recover RNN

- If $R_t \approx 0$, \tilde{h}_t only depends on x_t , resetting hidden states

Gated Recurrent Units (GRU)

UC San Diego

Reset gate: $r_t = \text{sigmoid}(W_{rx}x_t + W_{rh}h_{t-1} + b_r)$

Update gate: $z_t = \text{sigmoid}(W_{zx}x_t + W_{zh}h_{t-1} + b_z)$

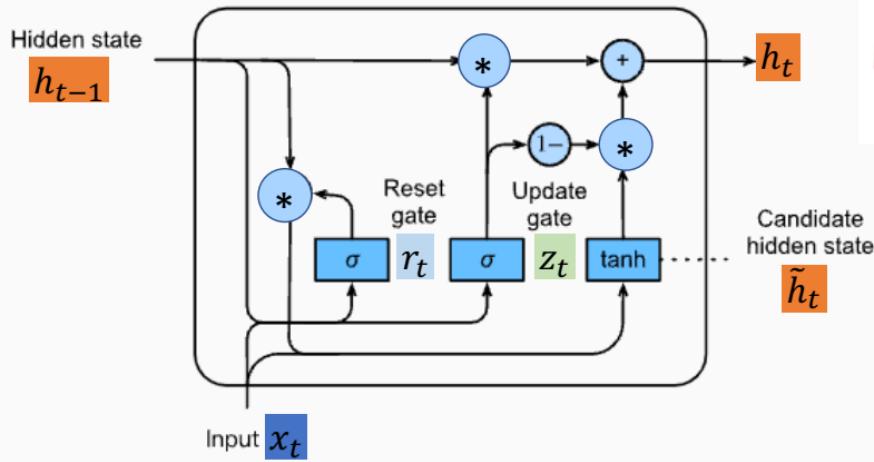
proposal $\tilde{h}_t = \tanh(W_{hx}x_t + W_{hh}(r_t * h_{t-1}) + b_h)$

Hidden state: $h_t = z_t * h_{t-1} + (1 - z_t) * \tilde{h}_t$
elementwise

- If $z_t \approx 1$, $h_t = h_{t-1}$; (help cope the vanishing gradient)
- If $z_t \approx 0$, $h_t = \tilde{h}_t$ (incorporate new info)

Gated Recurrent Units (GRU)

UC San Diego



$$\text{Reset gate: } r_t = \text{sigmoid}(W_{rx}x_t + W_{rh}h_{t-1} + b_r)$$

$$\text{Update gate: } z_t = \text{sigmoid}(W_{zx}x_t + W_{zh}h_{t-1} + b_z)$$

$$\text{Proposal: } \tilde{h}_t = \tanh(W_{hx}x_t + W_{hh}(r_t * h_{t-1}) + b_h)$$

$$\text{Hidden State: } h_t = z_t * h_{t-1} + (1 - z_t) * \tilde{h}_t$$

↑
elementwise

- Reset gates help capture short-term dependencies
- Update gates help capture long-term dependencies



FC layer with
activation function



Elementwise
operator



Copy

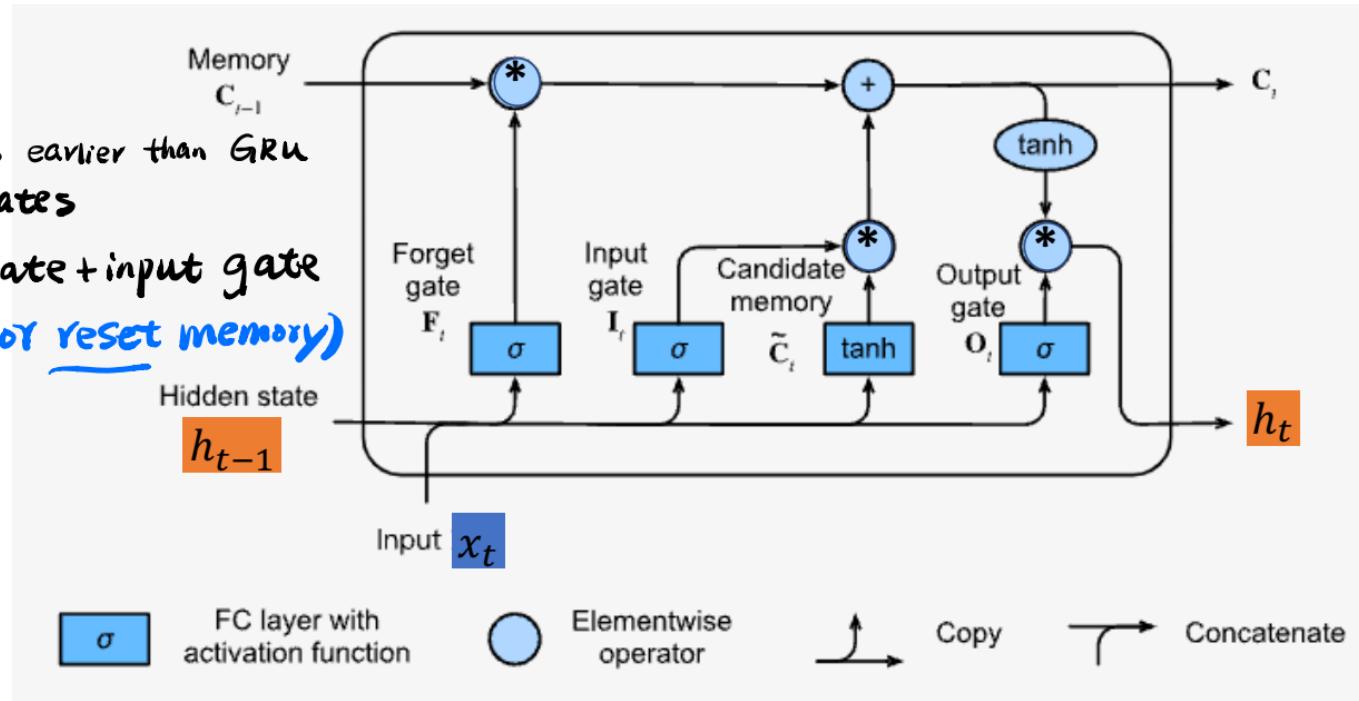


Concatenate

Long Short-Term Memory (LSTM)

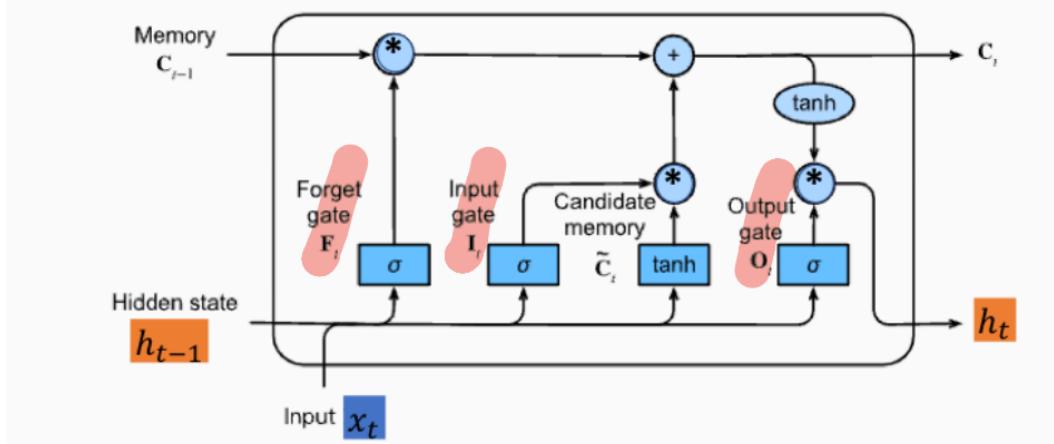
UC San Diego

- ~20 years earlier than GRU
- three-gates
- Forget gate + input gate
(retain or reset memory)
- Output gate



Long Short-Term Memory (LSTM)

UC San Diego



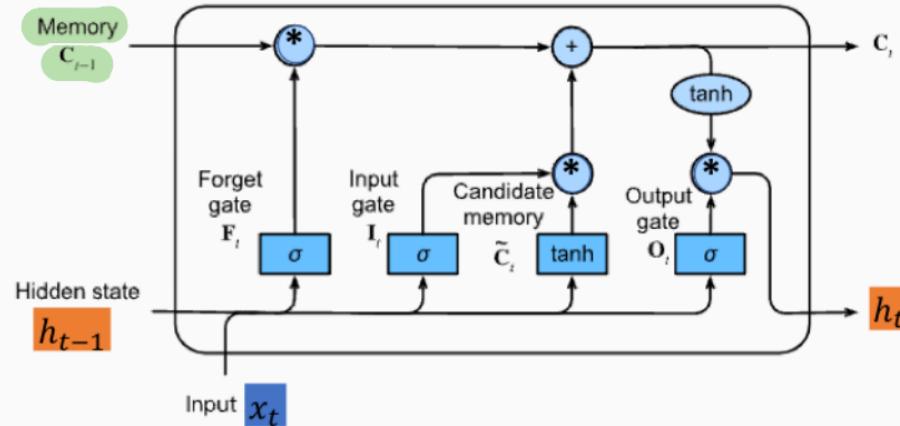
Input Gate: $I_t = \text{sigmoid}(W_{ix}x_t + W_{ih}h_{t-1} + b_i)$

Forget Gate: $F_t = \text{sigmoid}(W_{fx}x_t + W_{fh}h_{t-1} + b_f)$

Output Gate: $O_t = \text{sigmoid}(W_{ox}x_t + W_{oh}h_{t-1} + b_o)$

Long Short-Term Memory (LSTM)

UC San Diego



- I_t : how much we take x_t into account in computing \tilde{c}_t
- F_t : how much c_{t-1} is used

Proposal/Candidate. $\tilde{c}_t = \tanh(W_{cx}x_t + W_{ch}h_{t-1} + b_c)$

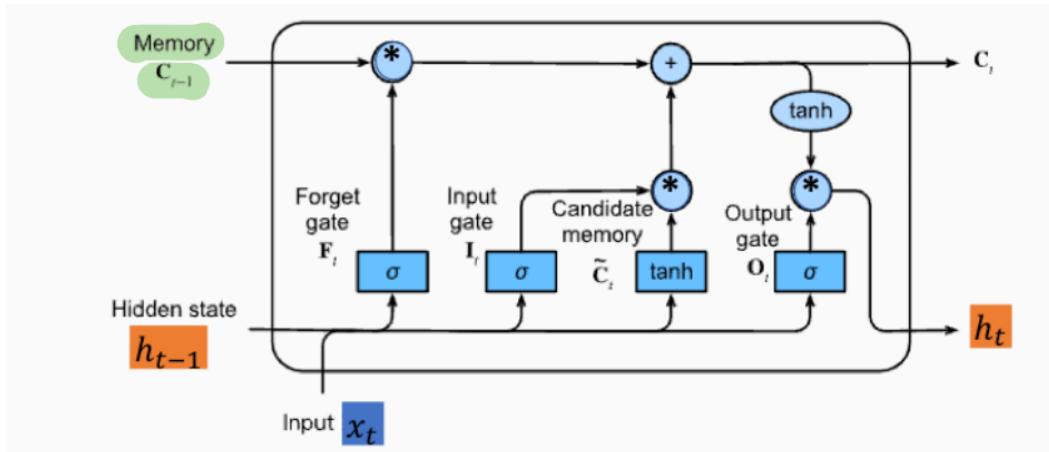
$c_t = F_t * c_{t-1} + I_t * \tilde{c}_t$, $*$ (element-wise multiplication)

Long short-term memory, <http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.676.4320&rep=rep1&type=pdf>

(GRU: $h_t = Z_t * h_{t-1} + (1 - Z_t) \tilde{h}_t$)

Long Short-Term Memory (LSTM)

UC San Diego



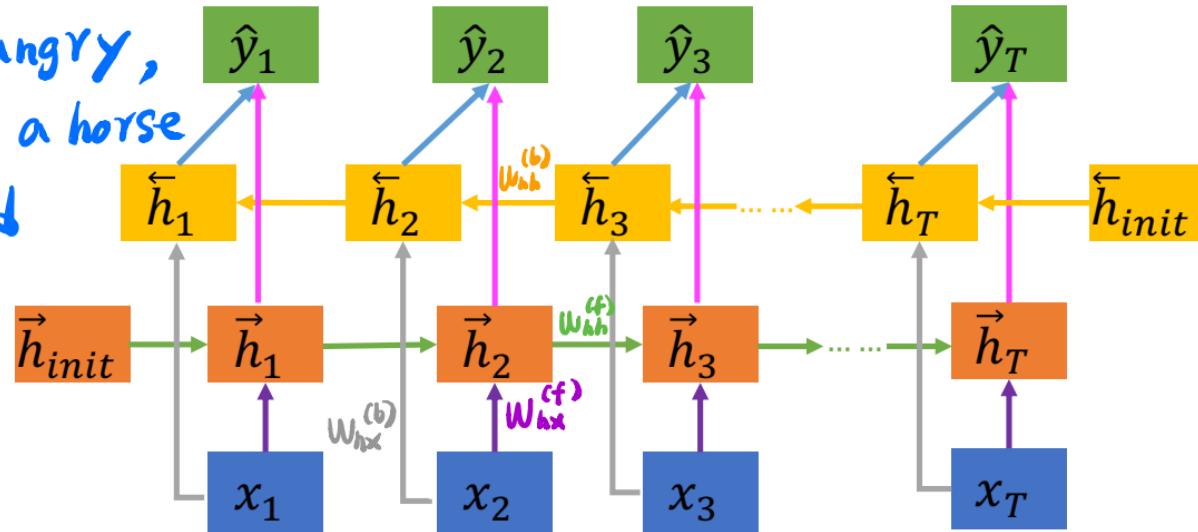
$$h_t = O_t * \tanh(c_t)$$

$$y_t = W_y h_t + b_y \quad (+ \text{activation func if needed})$$

Bidirectional RNN

UC San Diego

I am — hungry,
 - I can eat a horse
 - I just had
 dinner



$$\vec{h}_t = \sigma(W_{hh}^{(f)} \vec{h}_{t-1} + W_{hx}^{(f)} x_t + b_h^{(f)}) \quad x_t \in \mathbb{R}^{d \times n}, W_{hx} \in \mathbb{R}^{h \times d}, W_{hh} \in \mathbb{R}^{h \times h}$$

$$\overleftarrow{h}_t = \sigma(W_{hh}^{(b)} \overleftarrow{h}_{t+1} + W_{hx}^{(b)} x_t + b_h^{(b)}) \quad h_t = [\vec{h}_t, \overleftarrow{h}_t] \in \mathbb{R}^{2h \times n}$$

$$\hat{y}_t = \sigma(W_y h_t + b_y), \quad W_y \in \mathbb{R}^{q \times 2h}, b_y \in \mathbb{R}^{q \times 1}$$

Outline

UC San Diego

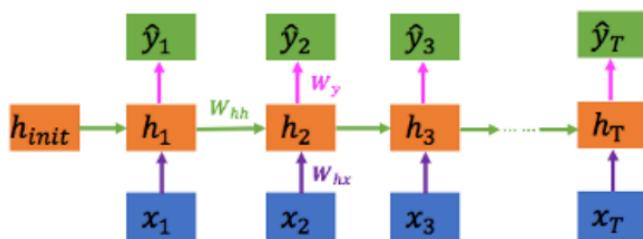
- Recurrent Neural Networks
- Modern Recurrent Neural Networks
- **Attention Mechanism and Transformer**

Motivation

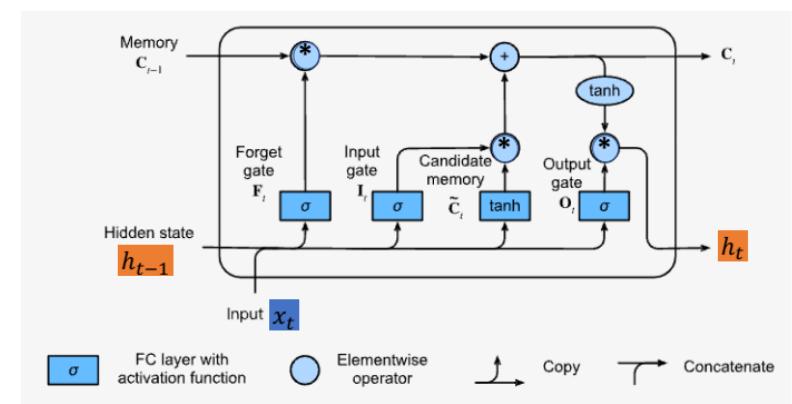
UC San Diego

Increased complexity,
sequential

RNN



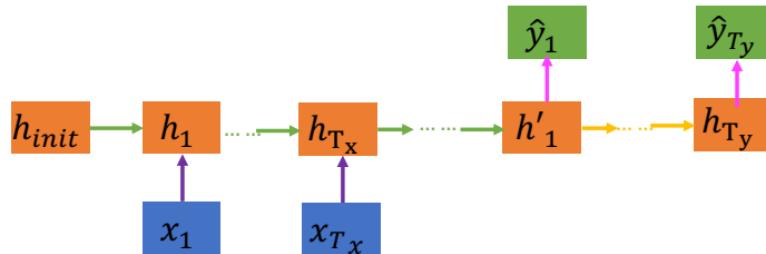
GRU



LSTM

Attention Mechanism

UC San Diego



The problem of long sequence in machine translation

DETECT LANGUAGE

ENGLISH

SPANISH

FRENCH



Scientists have been studying attention in the cognitive neuroscience field since the 19th century. In this chapter, we will begin by reviewing a popular framework explaining how attention is deployed in a visual scene. Inspired by the attention cues in this framework, we will design models that leverage such attention cues.



CHINESE (SIMPLIFIED)

ENGLISH

SPANISH



自 19 世纪以来，科学家们一直在研究认知神经科学领域的注意力。在本章中，我们将首先回顾一个流行的框架，解释如何在视觉场景中部署注意力。受此框架中的注意力线索的启发，我们将设计利用这些注意力线索的模型。



Attention Mechanism

UC San Diego



YouTube

physics informed machine learning



FILTERS

Top 6 Data Analytics Trends - Download Your Ebook Now
Learn To Simplify Your Complex Data That Will Provide Forecasting Business Opportunities.

Live Product Demos Cloud Data Warehousing Virtual Hands-On Labs Pricing Options

Ad <https://www.snowflake.com/online-ebook/data-science> VISIT SITE

Representing the solution of the PDE with a DNN

A Hands-on Introduction to Physics-informed Machine Learning
14K views • 10 months ago

nanohubtechtalks

2021.05.26 Ilias Bilionis, Atharva Hans, Purdue University Table of Contents below. This video is part of NCN's Hand

Deep Learning to Discover Coordinates for Dynamics: Autoencoders & Physics Informed Machine Learning
69K views • 8 months ago

Steve Brunton

Discovering physical laws and governing dynamical systems is often enabled by first learning a new coordinate sys...
the ...
4K CC 26:56

Learning a Discontinuous/Oscillatory Function in Physical & Fourier Domains

Fourier Spectrum

George Karniadakis - From PINNs to DeepONets
11K views • 1 year ago

Physics Informed Machine Learning

- **Query, Key, Value**

Example: Youtube video retrieval

- Query: "physics informed machine learning"
- Key: each video has a set of keys (video title, description, viewing count...)
- Values: video

Attention Mechanism

UC San Diego

Mathematically, suppose we have a query $q \in R^{d_q}$

and m key-value pairs $(k_1, v_1), \dots, (k_m, v_m)$

where $k_i \in R^k$ and each value $v_i \in R^{d_v}$.

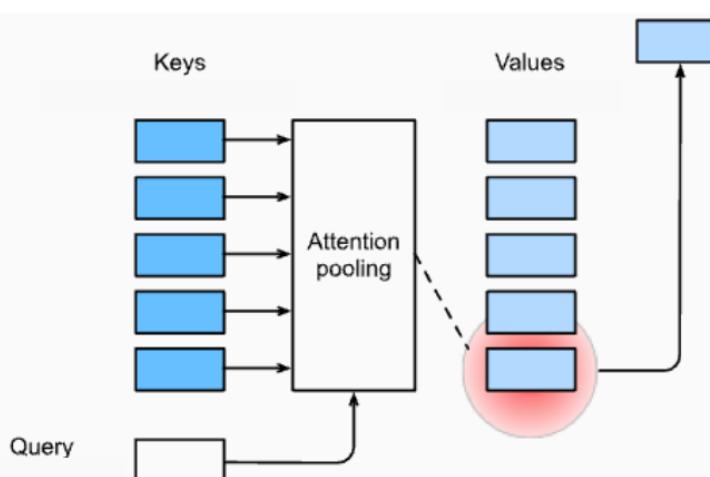
The output of the query can be written as:

Weight · amount of attention q should pay to item (k_i, v_i)

$$f(q, (k_1, v_1), \dots, (k_m, v_m)) = \sum_{i=1}^m \underline{\alpha(q, k_i)} v_i$$

$$\alpha(q, k_i) = \frac{\exp(\alpha(q, k_i))}{\sum_{j=1}^m \exp(\alpha(q, k_j))} \in R$$

α is called attention score function



Attention Mechanism

UC San Diego

Different attention score function:

- Query and key are vector of same length: $d_q = d_k$

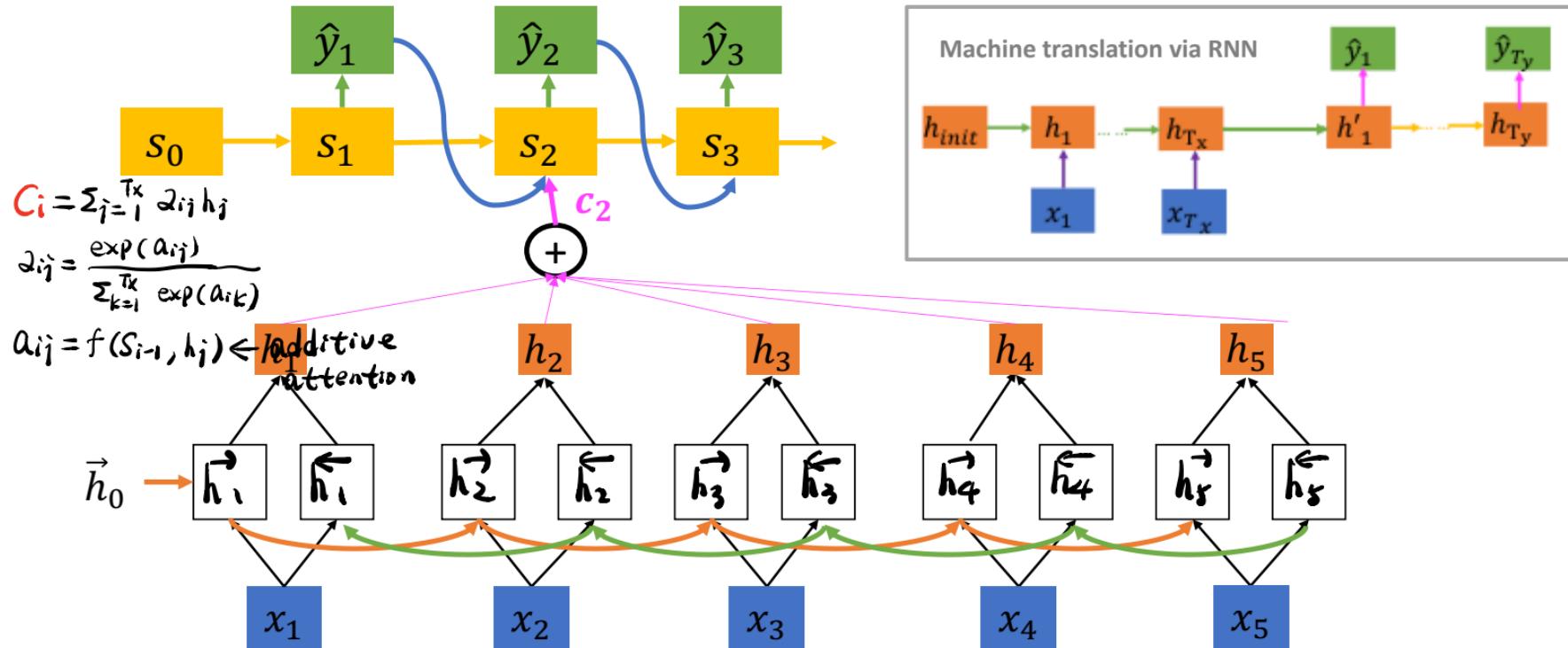
$$a(q, k_i) = \frac{q^T k_i}{\sqrt{d_k}} = \frac{k_i^T q}{\sqrt{d_k}} \quad (\text{Scaled dot-product attention})$$

- Query and key are vector of different length (additive attention)

$$a(q, k_i) = W_a^T \tanh(W_q q + W_k k_i) \in R$$
$$\begin{matrix} R^{da} \\ R^{daxde} \\ R^{daxde} \end{matrix}$$

RNN with Attention Mechanism

UC San Diego



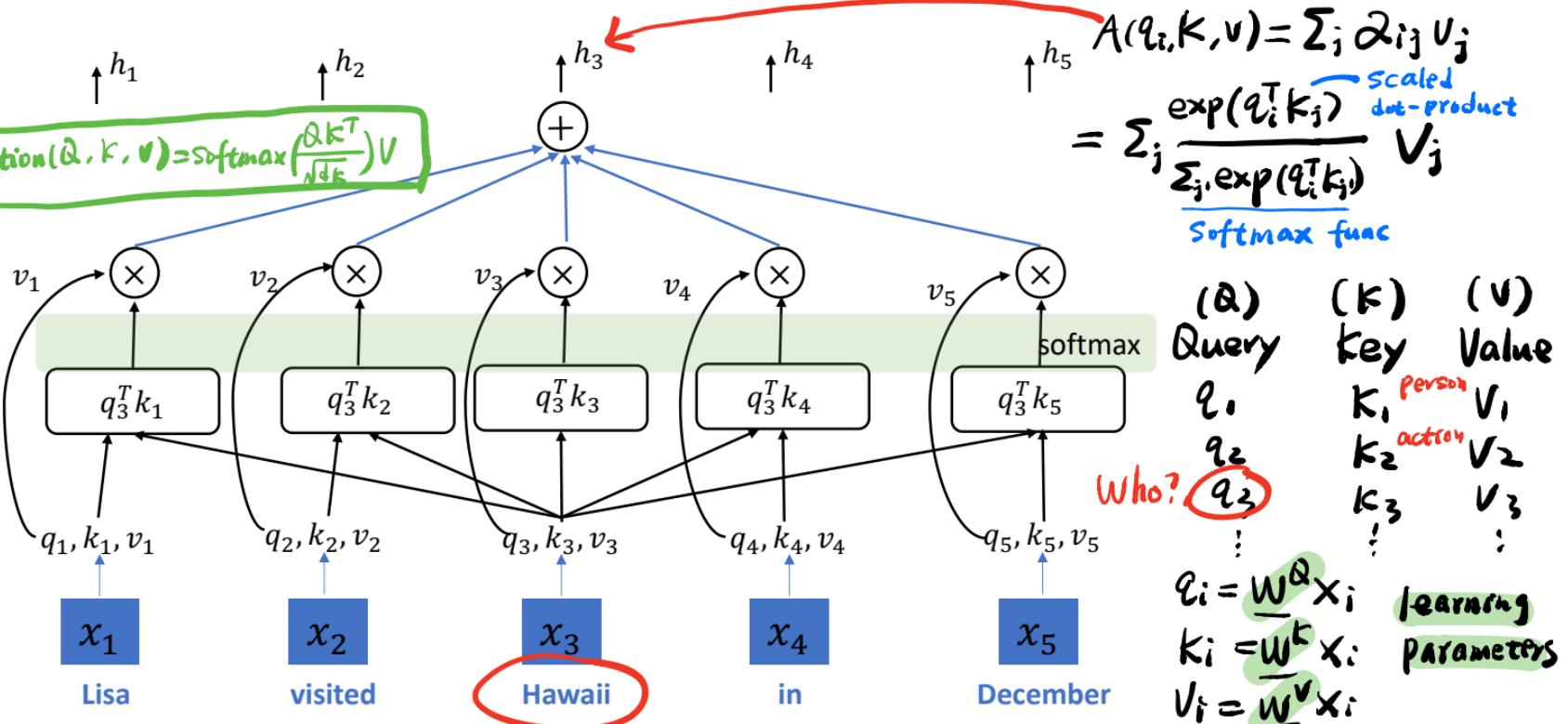
Transformer

UC San Diego

- Self Attention
- Multi-head Attention

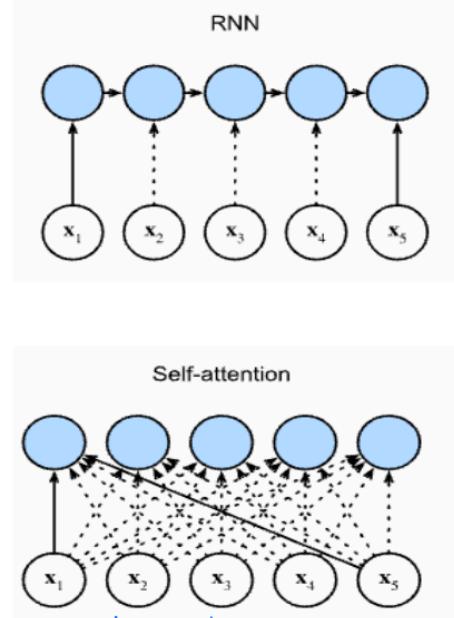
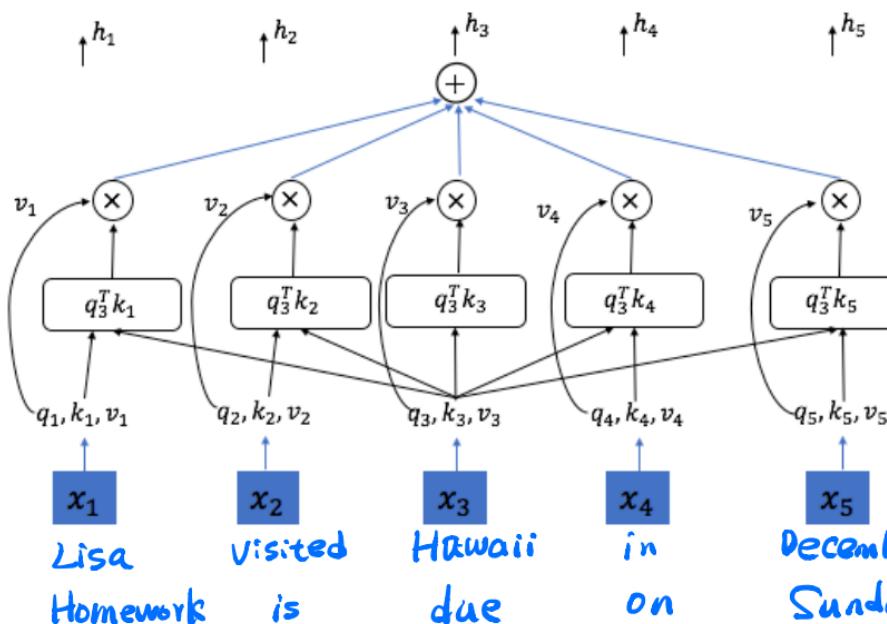
Transformer: Self Attention

UC San Diego



Transformer: Self Attention vs RNN

UC San Diego

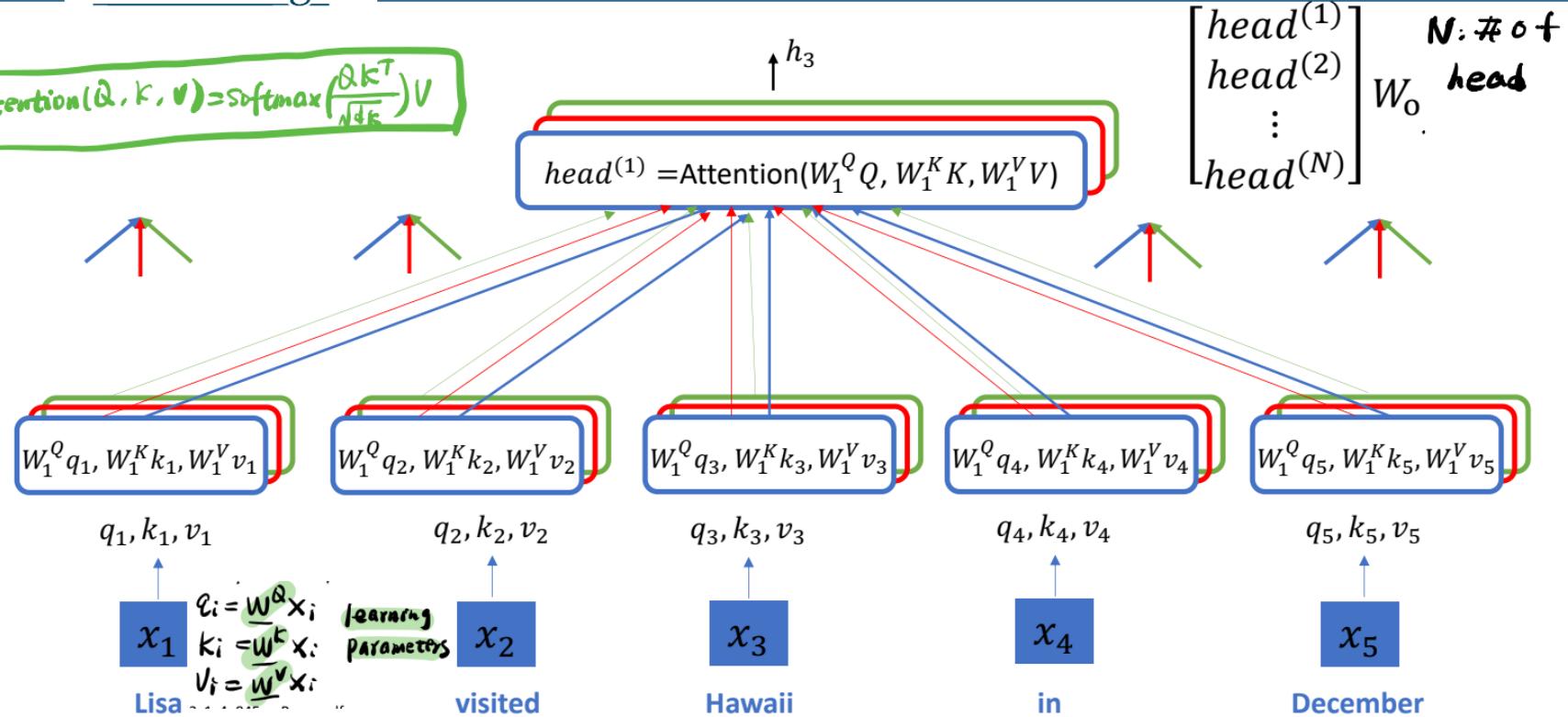


Attention Is All You Need, <https://proceedings.neurips.cc/paper/2017/file/3f5ee243547dee91fb053c1c4a845aa-Paper.pdf>

Transformer: Multi-Head Attention

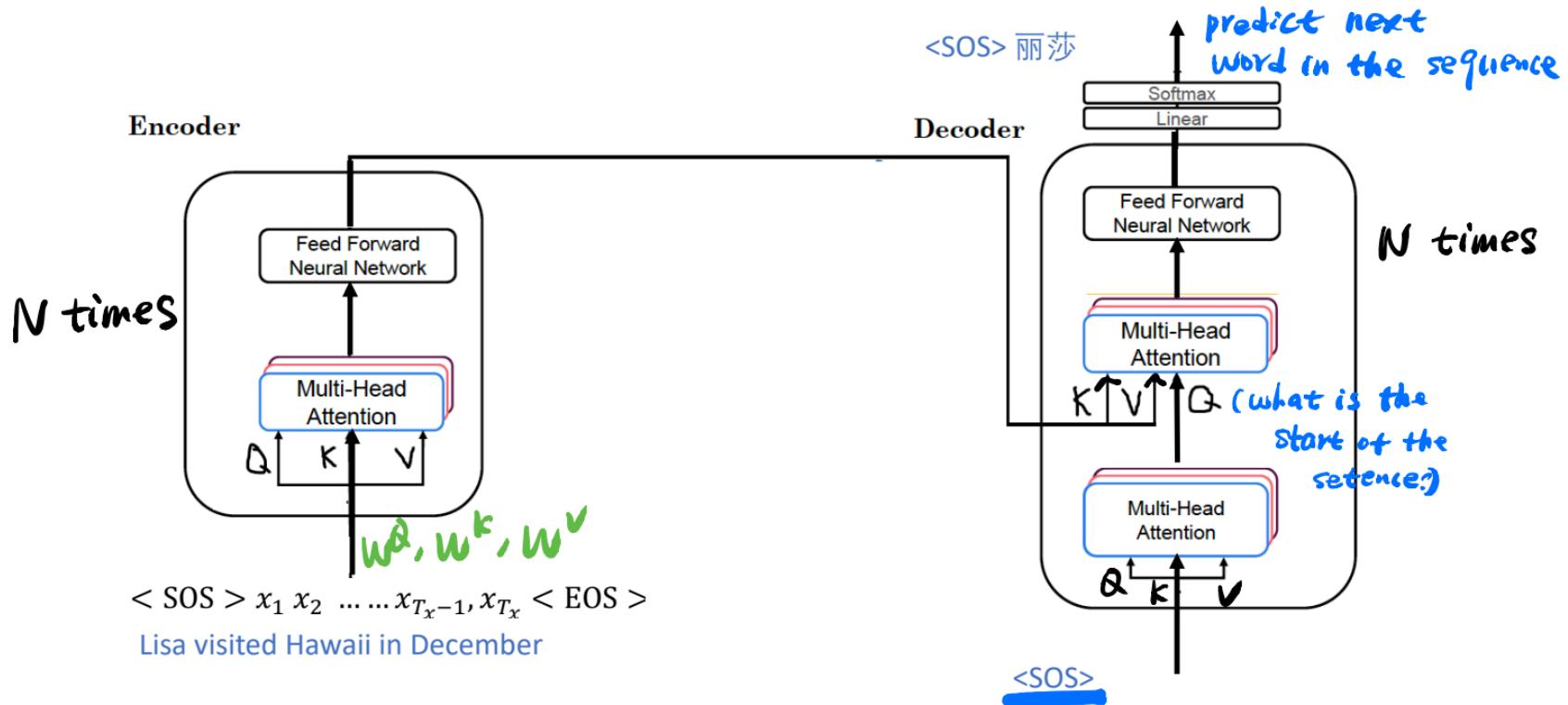
UC San Diego

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$



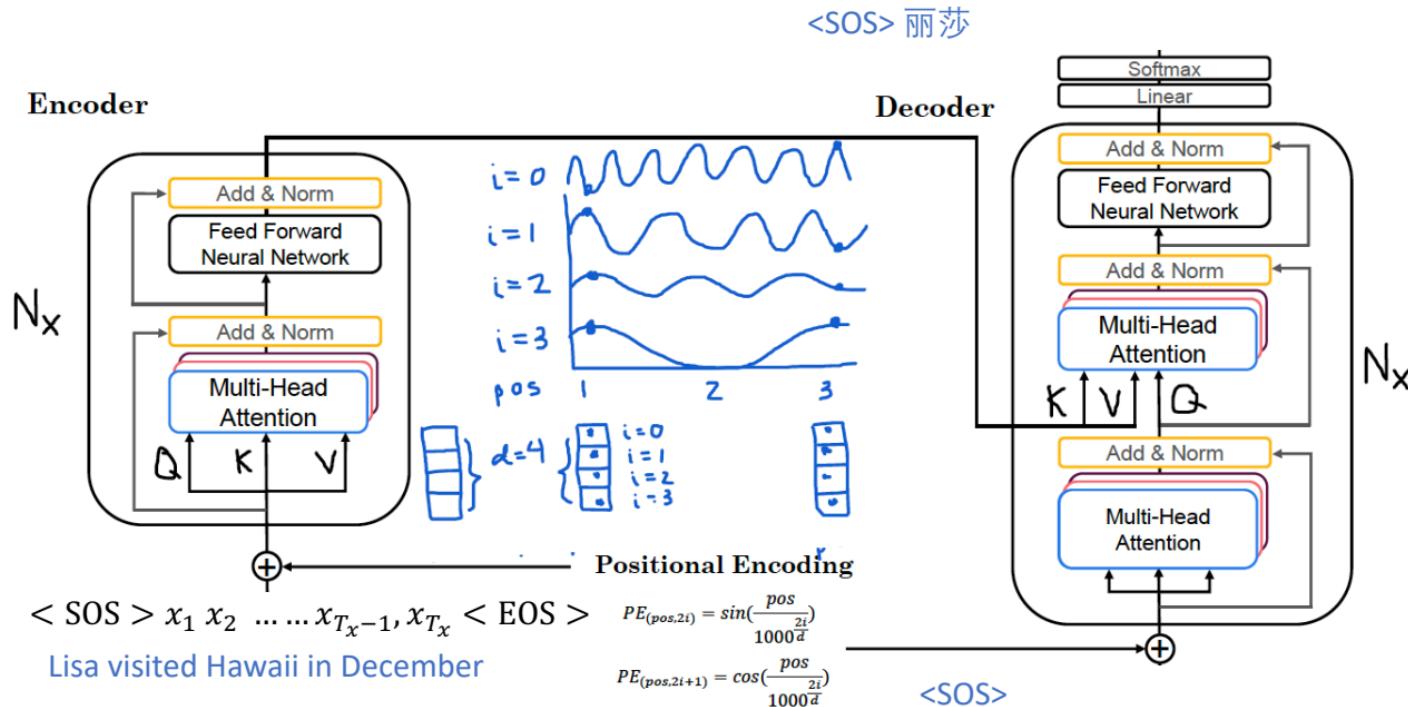
Transformer

UC San Diego



Transformer Details

UC San Diego



Transformer vs RNN

UC San Diego

Transformer

- * Convenient parallel computation
- * Can capture (very) long range dependency & not suffer from vanishing gradient issue
- * Work well when you have **BIG** data
(e.g. NLP tasks)

RNN / LSTM / GRU

- * better capture temporal sequence structure
- * RNN / LSTM / GRU can still work very well (or even better) in common tasks with moderate # of data and not hyper-long sequences