

# Recurrent neural networks

Vlad Gladkikh

IBS CMCM

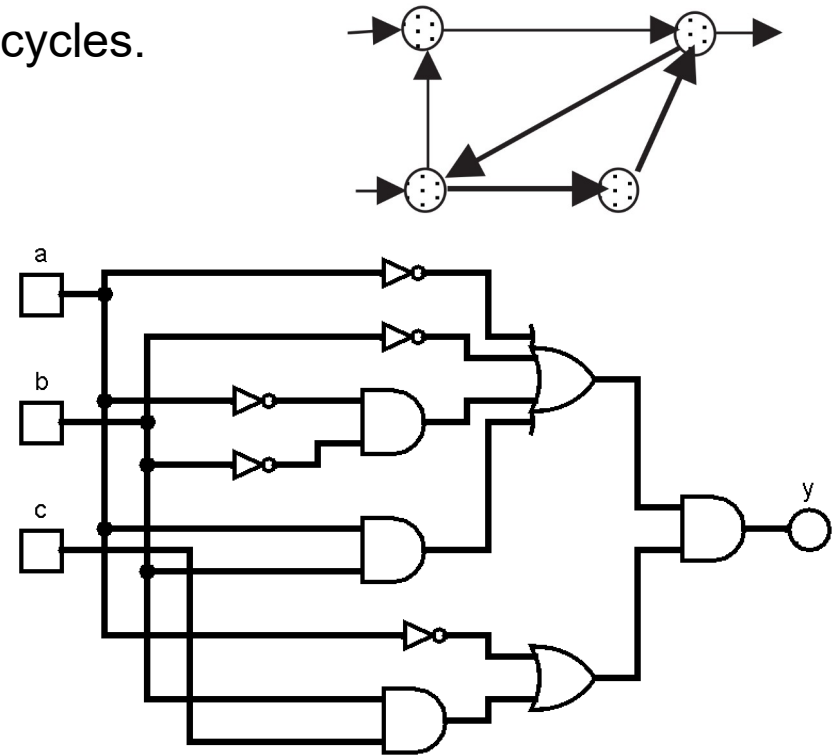
Recurrent neural networks are networks that contains cycles.

Similar concepts in digital electronics:

Two types of digital circuits

**Combinational circuits:** the outputs of the circuit depend only on the current inputs.

- do not need to know the history of past inputs
- do not require any memory elements



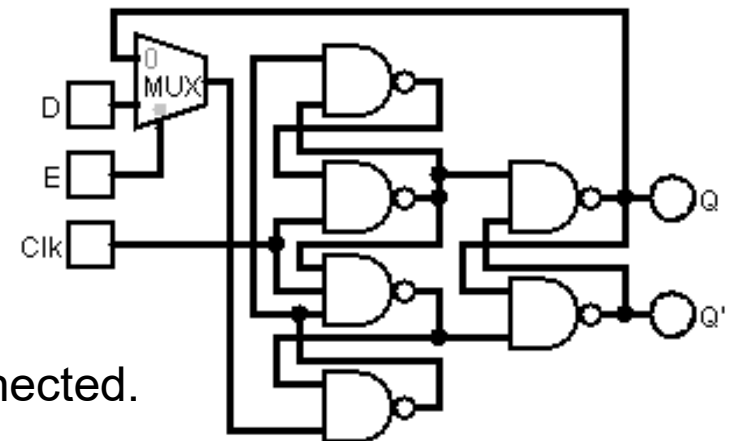
**Sequential circuits:** outputs depend on not only the current inputs, but also on all of the past inputs.

- must contain memory elements in order to remember the history of past input values

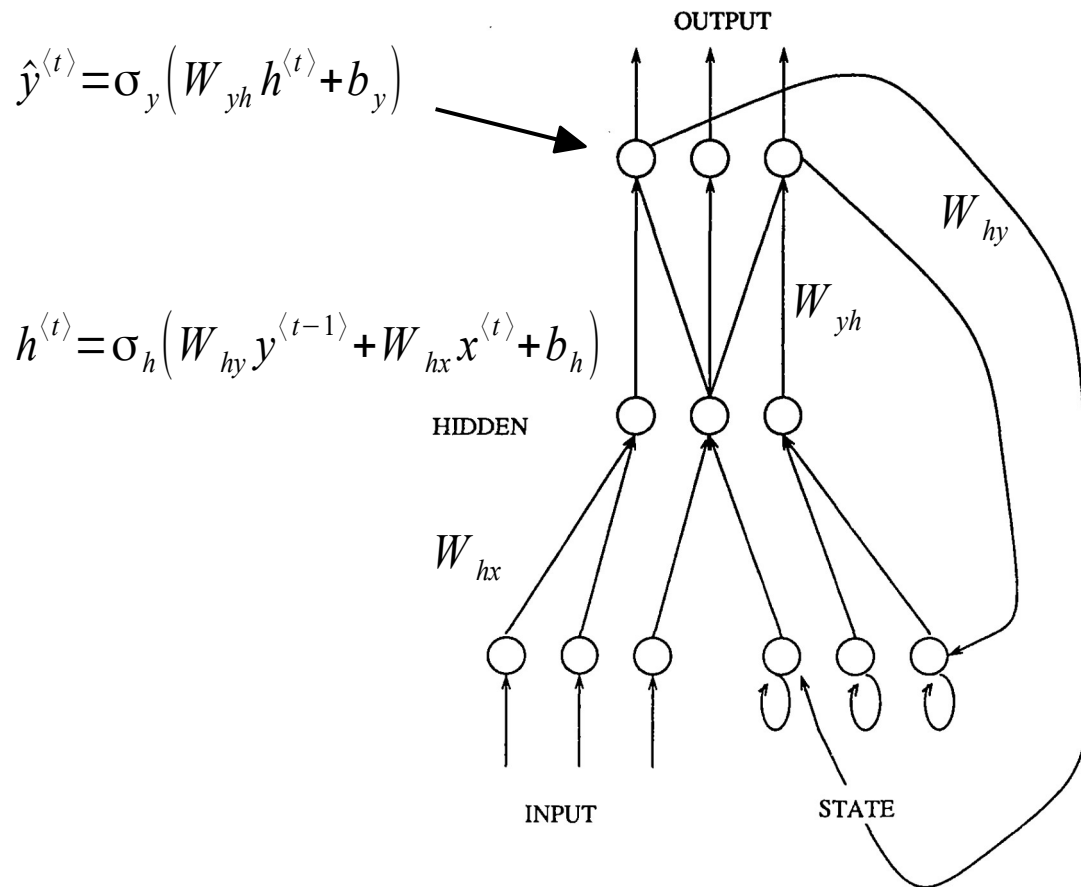
Both combinational and sequential circuits consist of the same building blocks – the AND, OR, and NOT gates, or just NAND gates.

What makes them different is in the way the gates are connected.

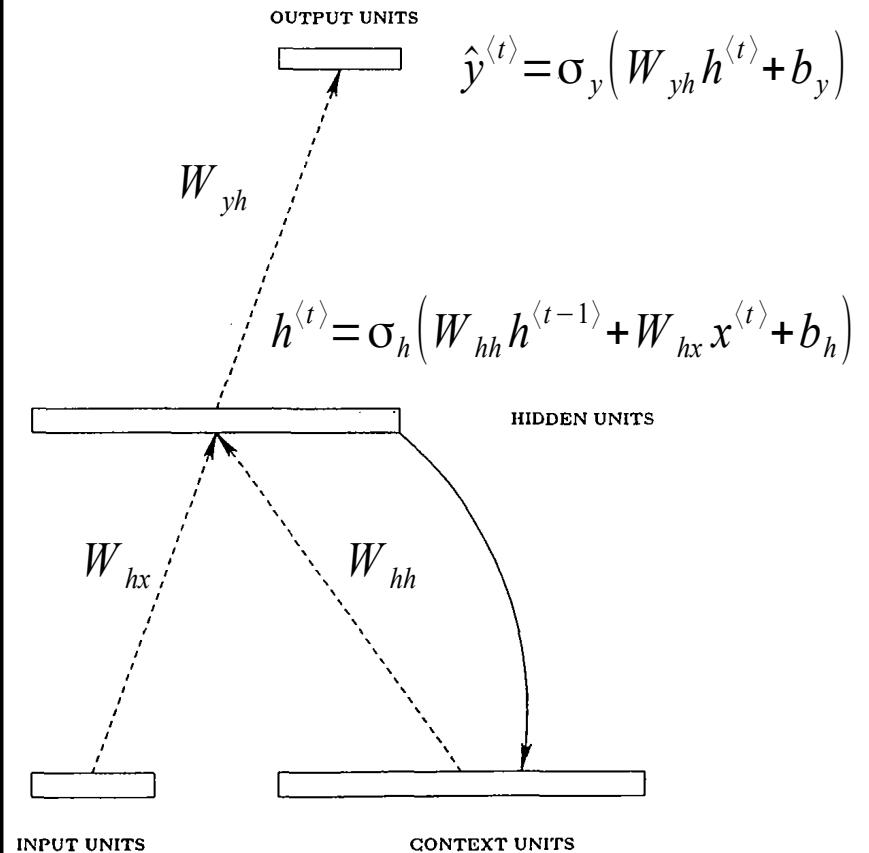
Sequential circuits have backward connections: output  $\rightarrow$  input



## Jordan network



## Elman network

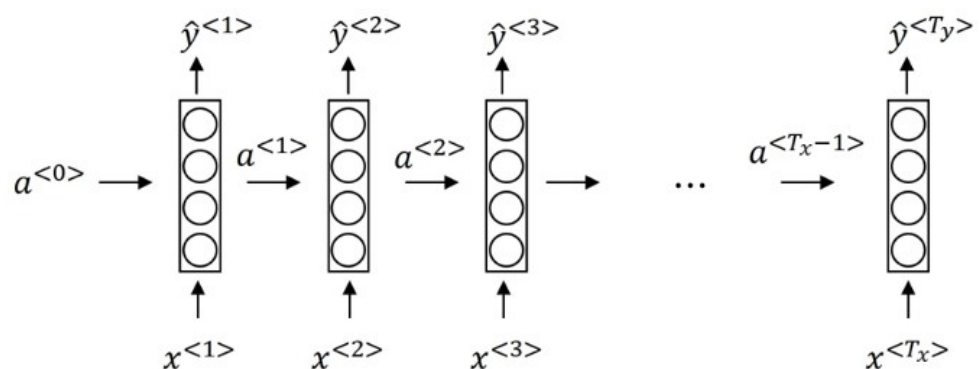


Dotted lines: trainable connections.

Elman, Jeffrey L. (1990). "Finding Structure in Time". Cognitive Science. 14 (2): 179–211. doi:10.1016/0364-0213(90)90002-E

<https://datascience.stackexchange.com/questions/82416/difference-between-jordan-elman-and-normal-rnn>

<https://ai.stackexchange.com/questions/8190/where-can-i-find-the-original-paper-that-introduced-rnns/>



$y^{(i)\langle t \rangle}$  – element  $t$  of output sample  $i$   
 $i = 1, \dots, \langle T_y \rangle$

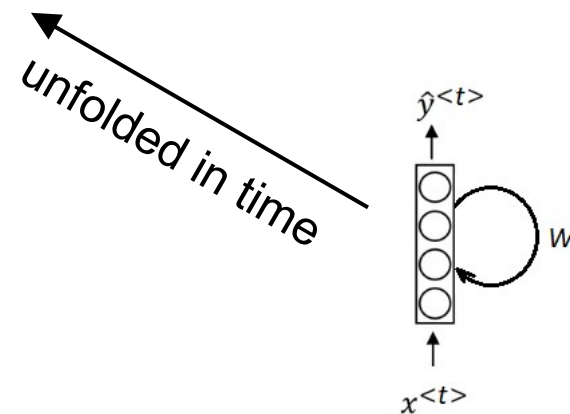
$x^{(i)\langle t \rangle}$  – element  $t$  of input sample  $i$   
 $i = 1, \dots, \langle T_x \rangle$

time

$$a^{\langle t \rangle} = g_a(W_{aa} a^{\langle t-1 \rangle} + W_{ax} x^{\langle t \rangle} + b_a) = g_a(W_a [a^{\langle t-1 \rangle}, x^{\langle t \rangle}] + b_a)$$

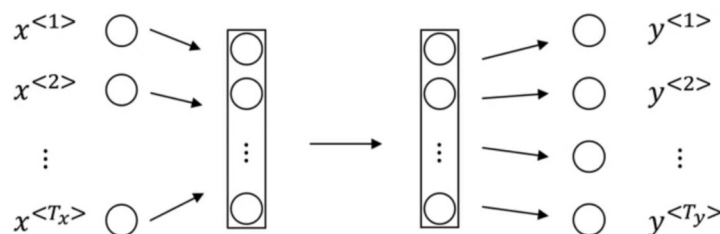
$$W_a = [W_{aa}, W_{ax}]$$

$$\hat{y}^{\langle t \rangle} = g_y(W_{ya} a^{\langle t \rangle} + b_y) = g_y(W_y a^{\langle t \rangle} + b_y)$$

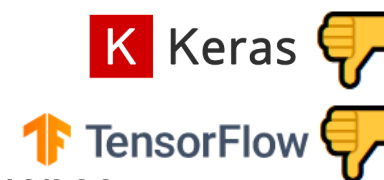


unfolded in time

Why not a feed-forward NN?

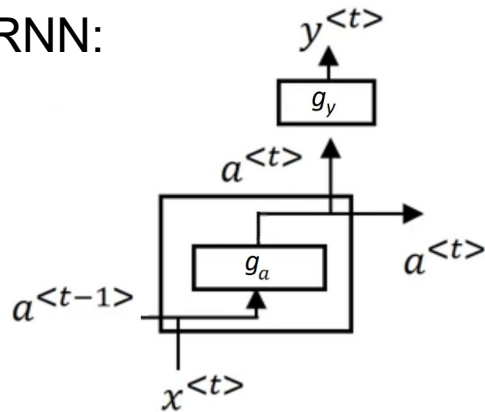


- Inputs and outputs can be different lengths in different examples
- No sharing of features learned across different positions in the sequence



$$\begin{bmatrix} 0 & 1 & 1 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 1 & 1 & 0 & 0 & 0 \end{bmatrix} \quad \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

RNN:



GRU and LSTM are RNN that can learn over long sequences.

<https://www.coursera.org/learn/nlp-sequence-models/>

## GRU - Gated Recurrent Unit

$$\tilde{c}^{<t>} = g_c(W_c[\Gamma_r \circ c^{<t-1>}, x^{<t>}] + b_c)$$

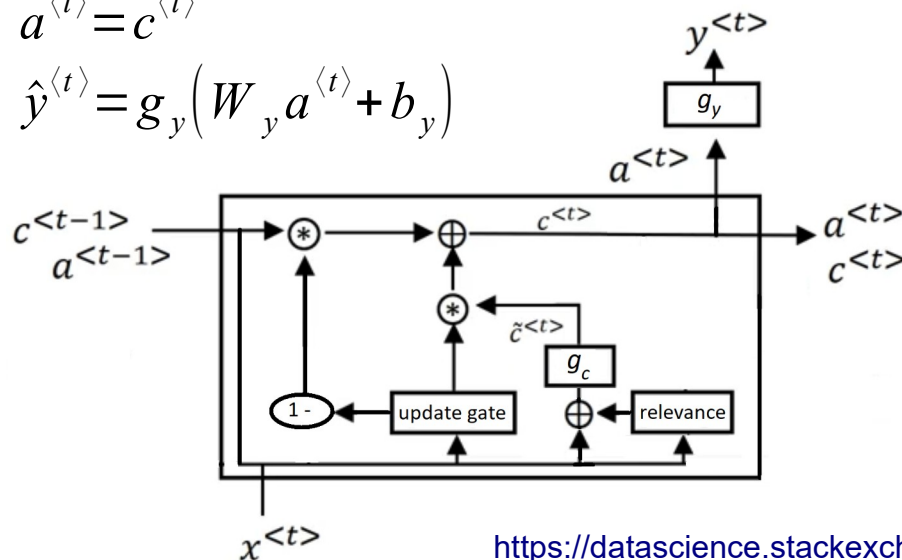
$$\Gamma_r = g_r(W_r[c^{<t-1>}, x^{<t>}] + b_r)$$

$$\Gamma_u = g_u(W_u[c^{<t-1>}, x^{<t>}] + b_u)$$

$$c^{<t>} = \Gamma_u \circ \tilde{c}^{<t>} + (1 - \Gamma_u) \circ c^{<t-1>}$$

$$a^{<t>} = c^{<t>}$$

$$\hat{y}^{<t>} = g_y(W_y a^{<t>} + b_y)$$



$$a^{<t>} = g_a(W_a[a^{<t-1>}, x^{<t>}] + b_a)$$

$$\hat{y}^{<t>} = g_y(W_y a^{<t>} + b_y)$$

$c$  - cell

$r$  - relevance

$u$  - update

$f$  - forget

$o$  - output

$\circ$  - element-wise product

Cell state: Long term memory

Hidden state: Working memory

## LSTM - Long Short-Term Memory

$$\tilde{c}^{<t>} = g_c(W_c[a^{<t-1>}, x^{<t>}] + b_c)$$

$$\Gamma_u = g_u(W_u[a^{<t-1>}, x^{<t>}] + b_u)$$

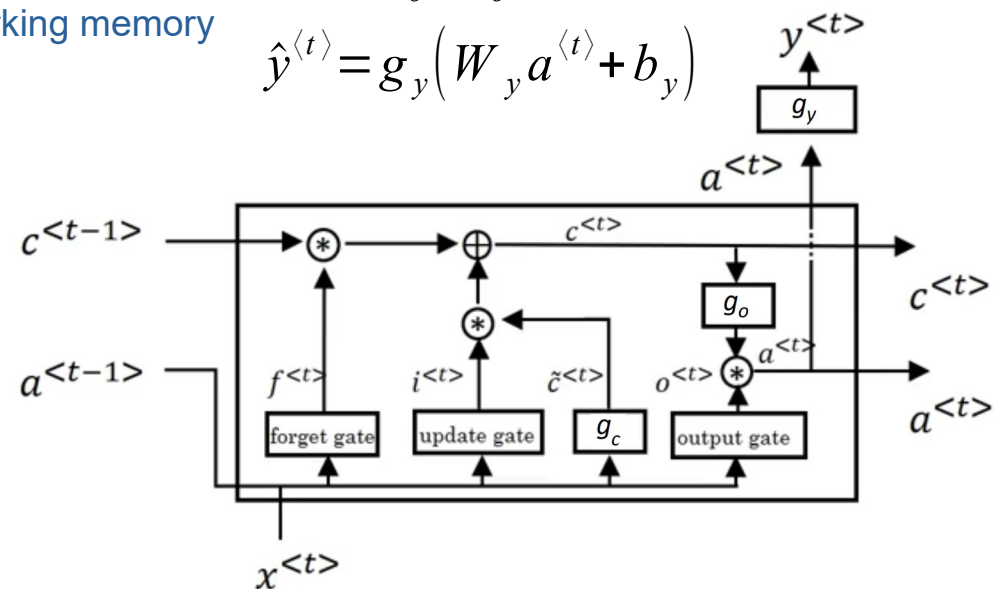
$$\Gamma_f = g_f(W_f[a^{<t-1>}, x^{<t>}] + b_f)$$

$$\Gamma_o = g_o(W_o[a^{<t-1>}, x^{<t>}] + b_o)$$

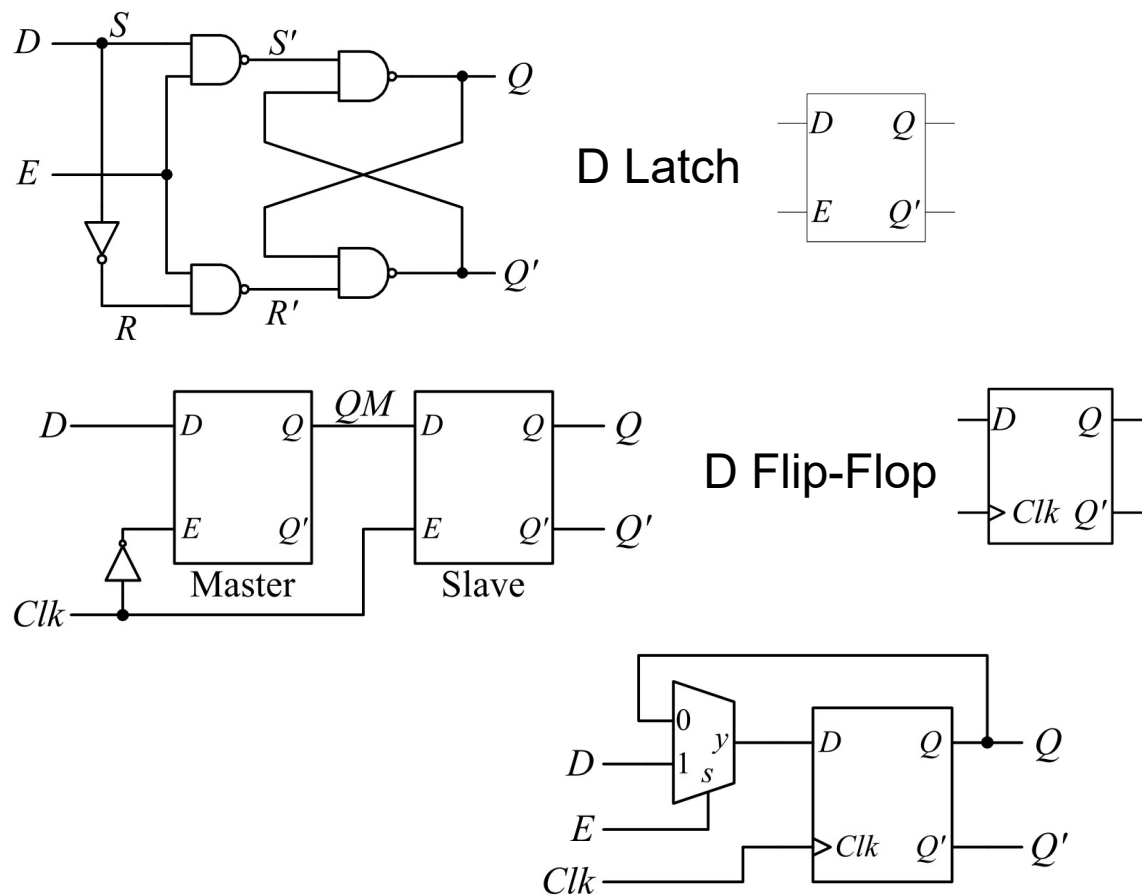
$$c^{<t>} = \Gamma_u \circ \tilde{c}^{<t>} + \Gamma_f \circ c^{<t-1>}$$

$$a^{<t>} = \Gamma_o \circ g_o(c^{<t>})$$

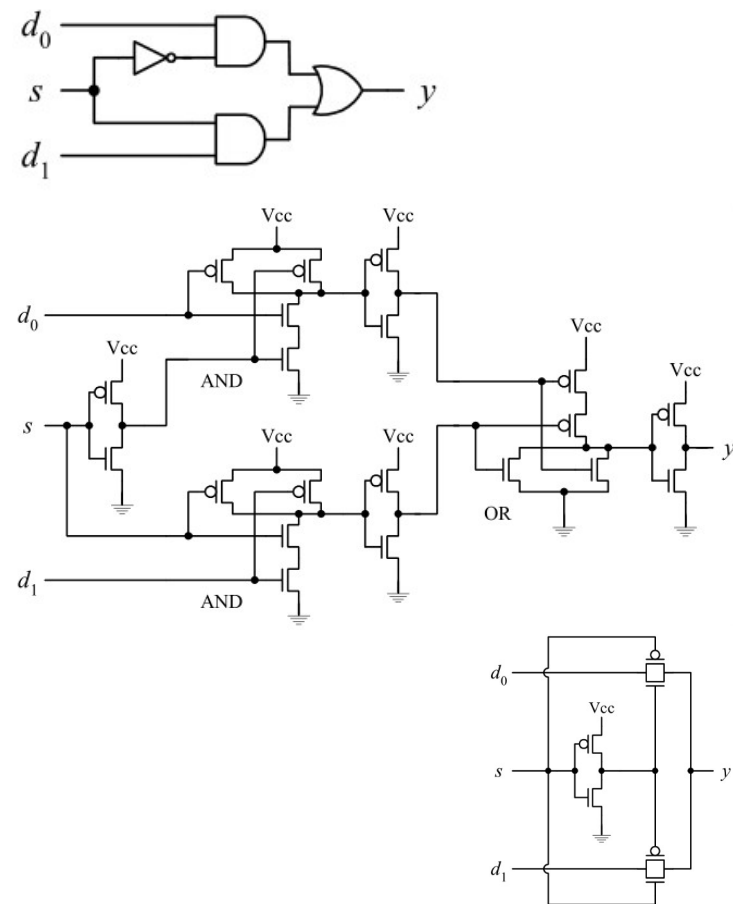
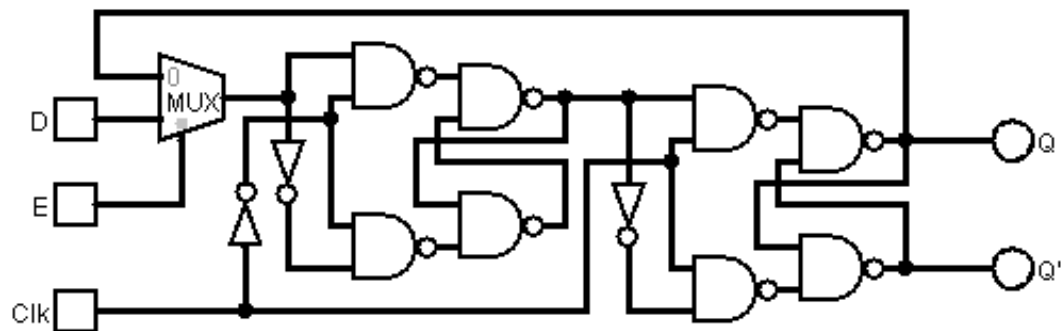
$$\hat{y}^{<t>} = g_y(W_y a^{<t>} + b_y)$$



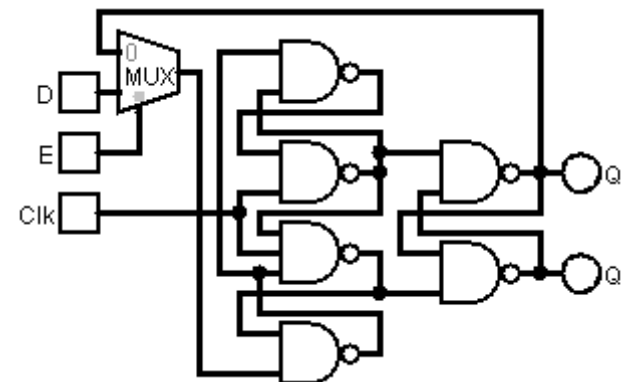
<https://datascience.stackexchange.com/questions/82808/difference-between-cell-state-and-hidden-state>



Naive implementation



A smarter implementation

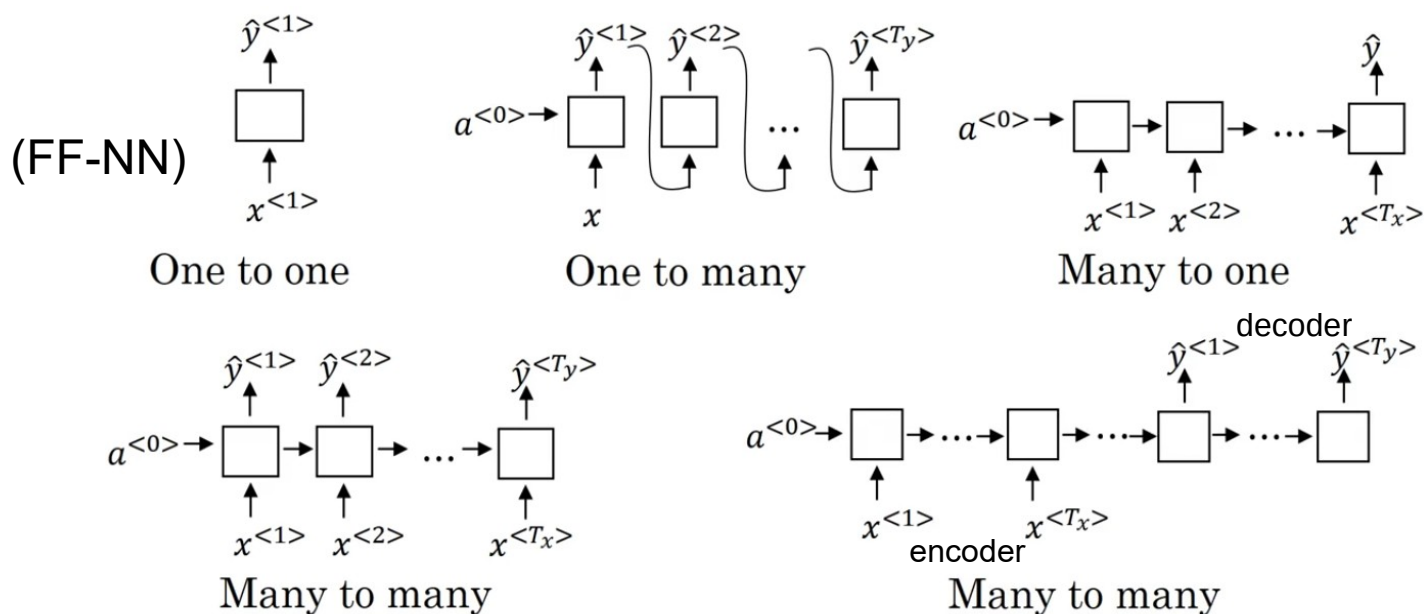


RNN can operate over sequences of vectors

<https://www.coursera.org/learn/nlp-sequence-models/>

Sequences can be in the input, the output, or both

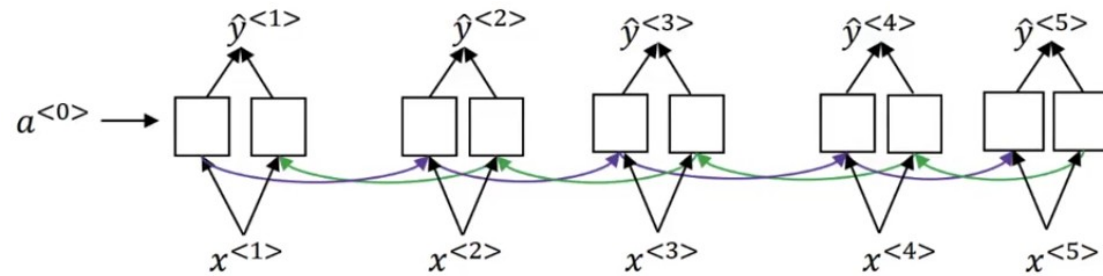
No constraints on the lengths sequences since the recurrent transformation is fixed and can be applied many times.



Implicit assumptions (wrong for some problems):

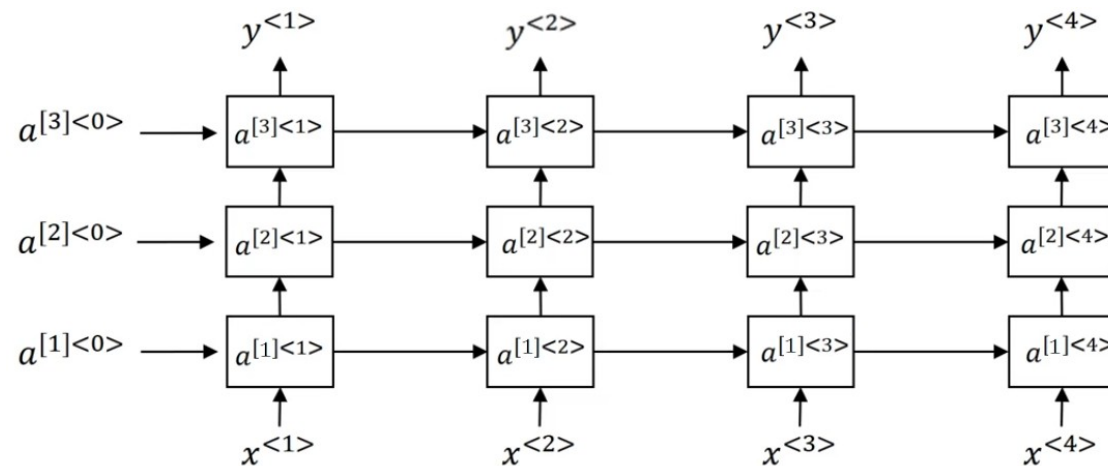
- 1) Only the preceding elements are relevant for the decision
- 2) The transformative relationship between the input, output and the passed-on state is the same for all elements in the chain, and doesn't change over time

## Bidirectional RNN



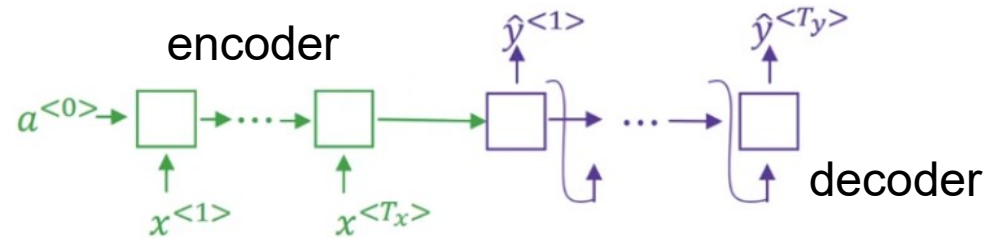
He said, "Teddy Roosevelt was a great President."  
He said, "Teddy bears are on sale!"

## Deep RNN





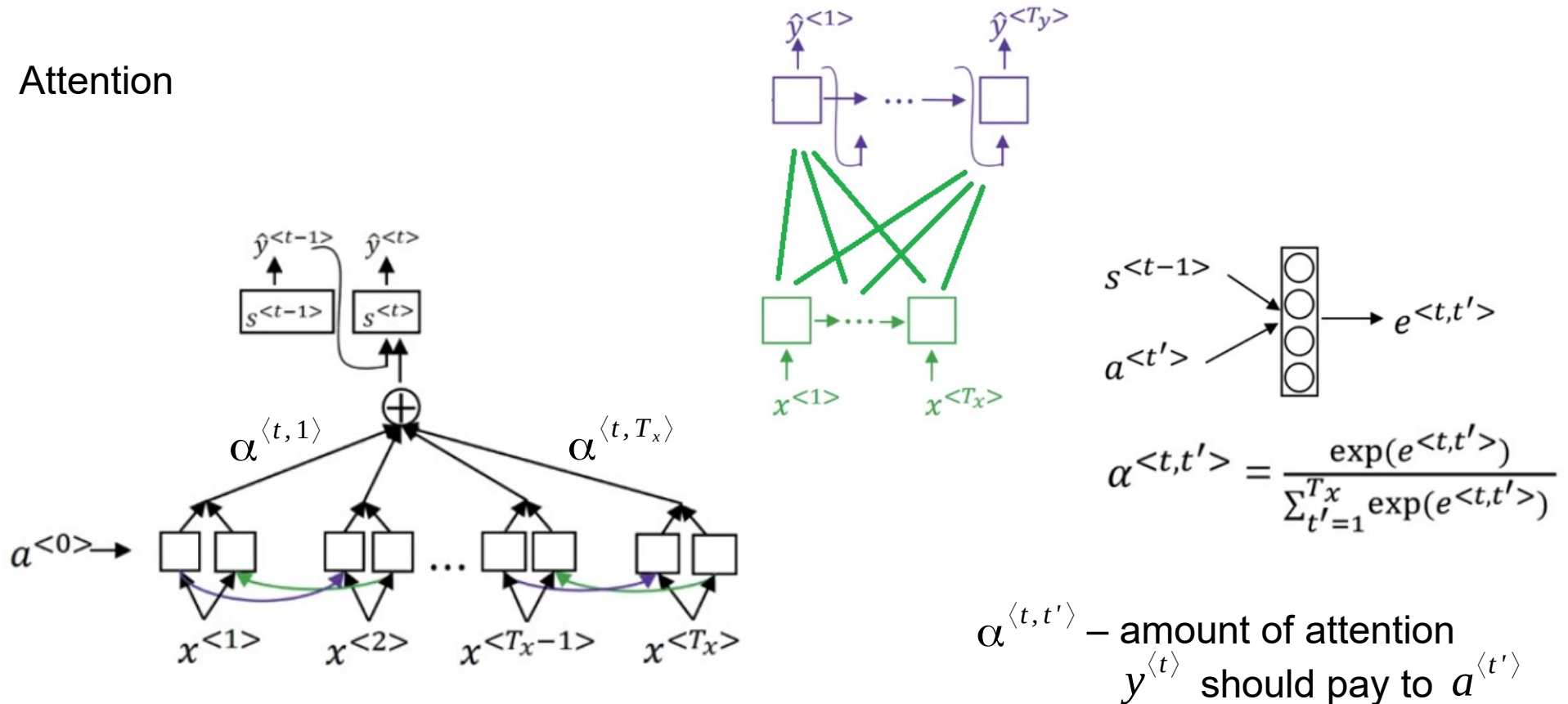
## Sequence to sequence



Jane s'est rendue en Afrique en septembre dernier, a apprécié la culture et a rencontré beaucoup de gens merveilleux; elle est revenue en parlant comment son voyage était merveilleux, et elle me tente d'y aller aussi.

Jane went to Africa last September, and enjoyed the culture and met many wonderful people; she came back raving about how wonderful her trip was, and is tempting me to go too.

## Attention



## Parts of sequential circuits

sequential circuit



### State memory

- keeps the complete history of inputs up to the current time.

The output of the state memory at one particular instance of time is called the **state** of the system at that time.

combinational circuit

### Output logic circuit

- generates the output of the whole circuit based on the current state, and may or may not be dependent on the current inputs.

combinational circuit

### Next-state logic circuit

- responsible for determining what next state to go to based on the current state that the system is in (i.e., the past inputs), and the current inputs

---

A sequential circuit operates by transitioning from one state to the next, generating different output signals.

# Word representation

<https://www.coursera.org/learn/nlp-sequence-models/>

$V = [a, aaron, \dots, zulu, <UNK>]$

## 1-hot representation

Man (5391)	Woman (9853)	King (4914)	Queen (7157)	Apple (456)	Orange (6257)
$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix}$

No inherent relationship to one another

Similar words are likely to have similar vectors

## Featurized representation: word embedding

	Man (5391)	Woman (9853)	King (4914)	Queen (7157)	Apple (456)	Orange (6257)
Gender	-1	1	-0.95	0.97	0.00	0.01
Royal	0.01	0.62	<u>0.93</u>	<u>0.95</u>	-0.01	0.00
Age	0.03	0.02	0.7	0.69	0.03	-0.02
Food	0.04	0.01	0.02	0.01	0.95	0.97
size	$\vdots$	$\vdots$				
cost						
alive						
verb						

$e_{5391}$        $e_{9853}$

I want a glass of orange juice.  
 I want a glass of apple juice.

Andrew Ng

## Language modeling

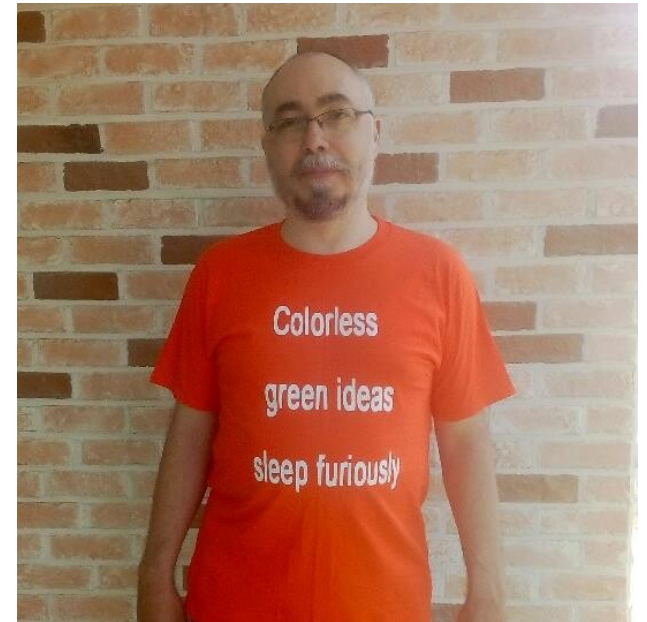
Given a sentence, what is the probability of that sentence?

E.g.

$$P(\text{The apple and pair salad}) = 3.2 \times 10^{-13}$$

$$P(\text{The apple and pear salad}) = 5.7 \times 10^{-10}$$

$$\frac{P(\text{Colorless green ideas sleep furiously.})}{P(\text{Furiously sleep ideas green colorless.})} = 2 \times 10^5$$



Tomas Mikolov, Wen-tau Yih, Geoffrey Zweig. Linguistic Regularities in Continuous Space Word Representations. Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, 2013, 746–751, <https://www.aclweb.org/anthology/N13-1090>

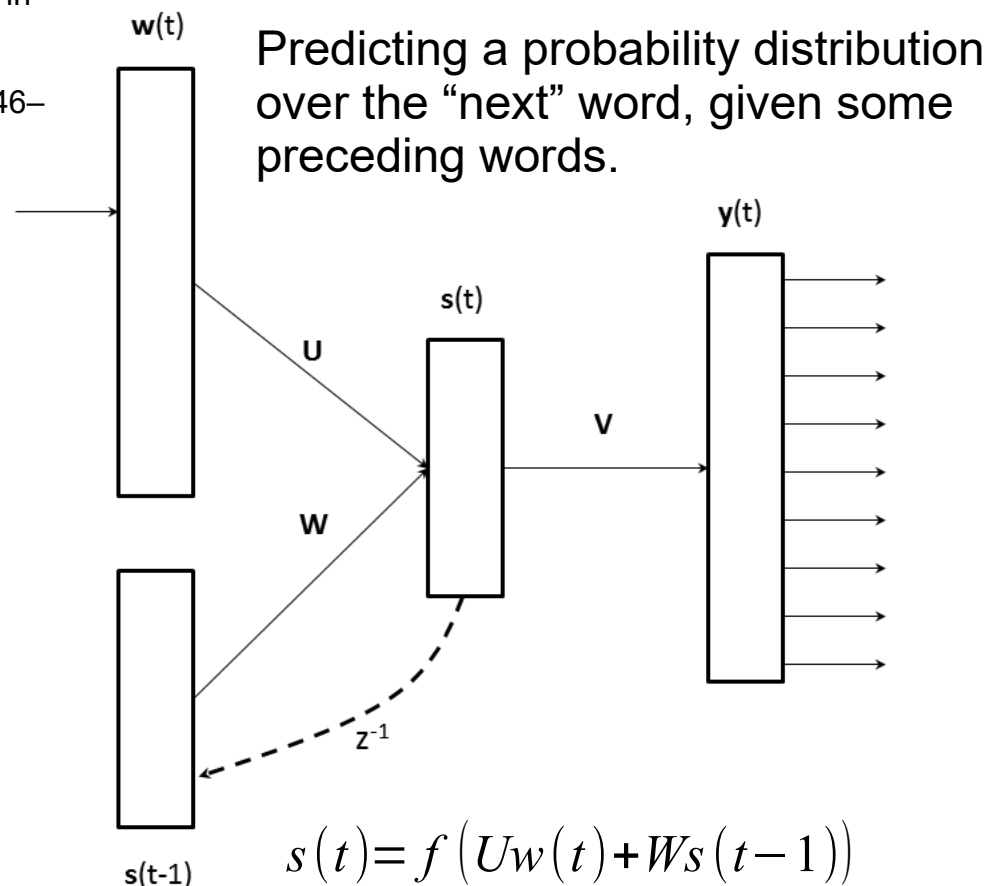
By training a neural network language model, one obtains not just the model itself, but also the learned word representations, which may be used for other, potentially unrelated, tasks.

$w(t)$  – input word at time  $t$  encoded using 1-of- $N$  coding

$y(t)$  – a probability distribution over words

$s(t)$  maintains a representation of the sentence history

$w(t), y(t)$  have dimensionality of the vocabulary

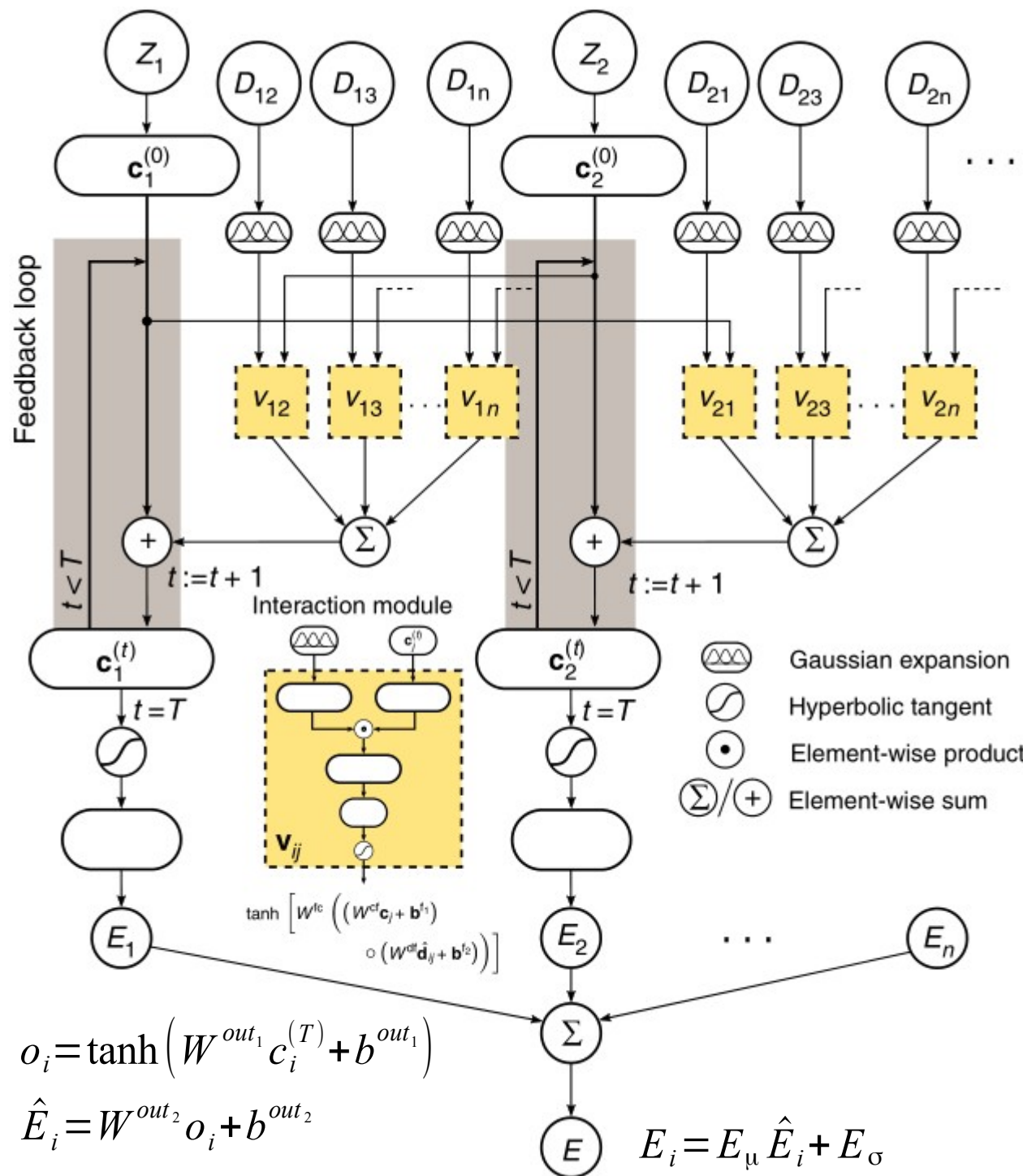


$$s(t) = f(Uw(t) + Ws(t-1))$$

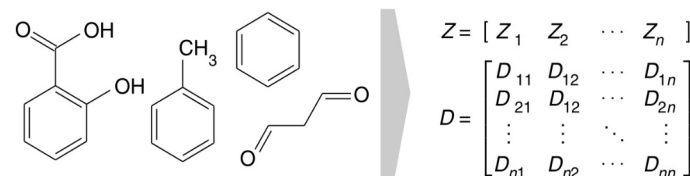
$$y(t) = g(Vs(t))$$

$$f(z) = \frac{1}{1 + e^{-z}} \quad g(z_m) = \frac{e^{z_m}}{\sum_k e^{z_k}}$$

The word representations are found in the columns of  $U$ , with each column representing a word.



K. T. Schütt et al. Quantum-Chemical Insights from Deep Tensor Neural Networks, Nat. Commun. 8, 13890, 2017.  
<https://www.nature.com/articles/ncomms13890>



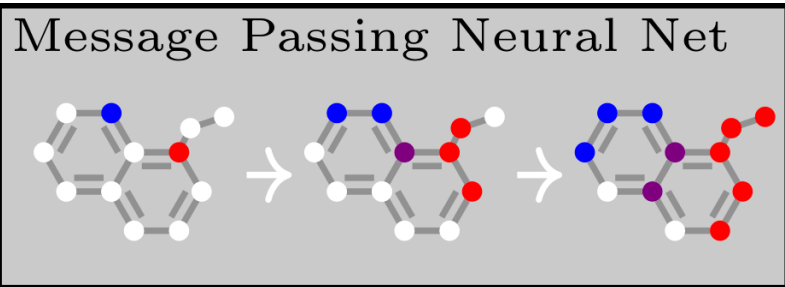
The distances are expanded in a Gaussian basis, yielding a feature vector  $\hat{d}_{ij} \in \mathbb{R}^G$

Each atom  $i$  is represented by a coefficient vector  $c \in \mathbb{R}^B$

$B$  : # of basis functions (features)

This atomic expansion is repeatedly refined by pairwise interactions with the surrounding atoms

$$c_i^{(t+1)} = c_i^{(t)} + \sum_{j \neq i} v_{ij}$$



Learning features from (molecular) graphs

$x_v$  node features

$e_{vw}$  edge features

$h_v^t$  hidden states for vertices

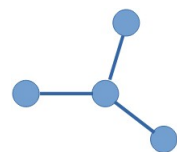
$h_{vw}^t$  hidden states for edges

Gilmer et al Neural Message Passing for Quantum Chemistry <https://arxiv.org/abs/1704.01212>

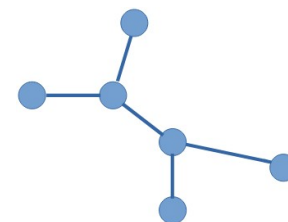
Two phases:

1) Message passing phase

$T$  time steps



message functions



$$m_v^{t+1} = \sum_{w \in N(v)} M_v^t(h_v^t, h_w^t, e_{vw})$$

$$m_{vw}^{t+1} = \sum_{e \in N(e)} M_{vw}^t(h_{vw}^t, e, h_v^t, h_w^t)$$

$$h_v^{t+1} = U_v^t(h_v^t, m_v^{t+1})$$

$$h_{vw}^{t+1} = U_{vw}^t(h_{vw}^t, m_{vw}^{t+1})$$

update functions

2) Readout phase

$$\hat{y} = R(\{h_v^T, h_{vw}^T | v, w \in G\})$$

readout function, must be invariant to permutations of the node states

E.g. DTNN  $M_v^t = \tanh(W^{fc}((W^{cf} h_v^t + b_1) \odot (W^{df} e_{vw} + b_2)))$

$$U_v^t = h_v^t + m_v^{t+1}$$

$$R = \sum_v NN(h_v^T)$$



Swanson et al. Deep learning for automated classification and characterization of amorphous materials. Soft Matter, 2020,16, 435-446  
<https://doi.org/10.1039/C9SM01903K>

