

Sequence Model

Hendrik Santoso Sugiarto

Capaian Pembelajaran

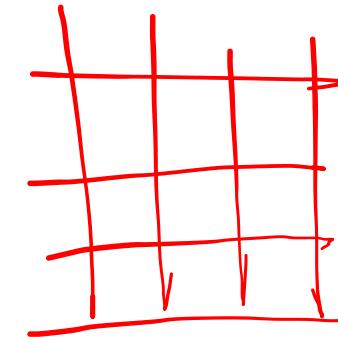
- Recurrent Neural Network
 - Mengerti dan mengaplikasikan RNN
- LSTM & RNN ~~GRU~~
 - Mengerti kelemahan RNN
 - Mengeerti dan menggunakan LSTM & GRU untuk menyelesaikan kelemahan RNN

RNN



Recurrent Neural Network

- CNN efektif dalam memproses grid data
- RNN efektif dalam memproses sequence data

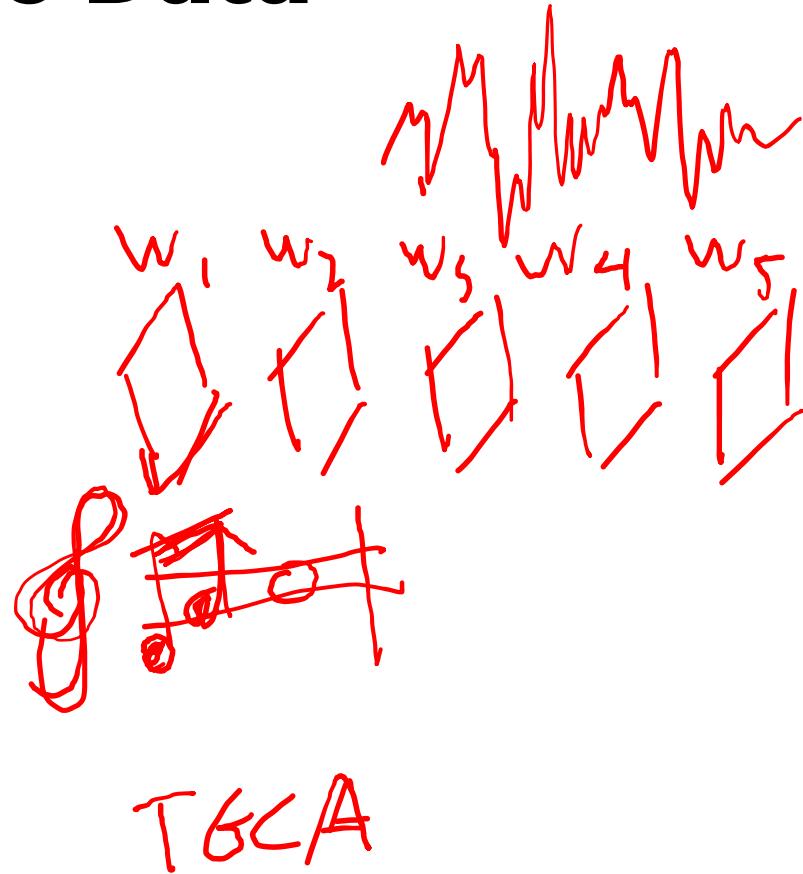


Motivasi

- Misalkan ada Neural Network yang memiliki input/output berikut
- Input: I am studying deep learning
- Output: Saya sedang belajar pembelajaran yang mendalam
- Contoh tersebut menunjukkan bahwa:
 - Konteks dan pola urutan sebelum-sesudah (sequence) sangat penting
 - Terdapat jenis data dengan input maupun output berbentuk sequence
 - Sulit untuk membuat fixed context window (mungkin ada kalimat yang lebih panjang dari sebelumnya)

Sequence Data

- Time-Series
- Text
- Video
- Music
- DNA



Basis Neurosains

MicroNetwork Motifs

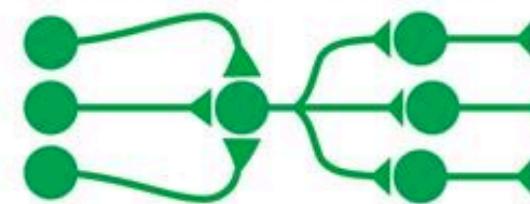
A. Feedforward excitation



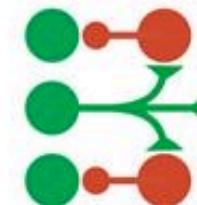
B. Feedforward inhibition



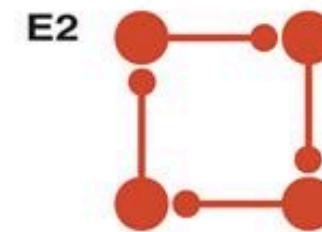
C. Convergence/divergence



D. Lateral inhibition



E. Feedback/Recurrent inhibition



F. Feedback/Recurrent excitation

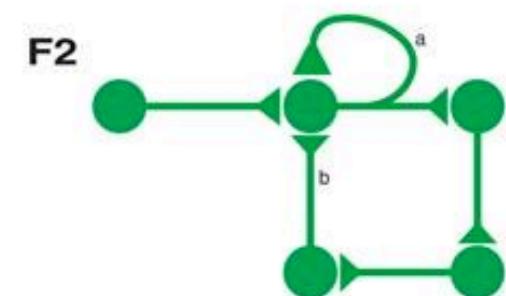
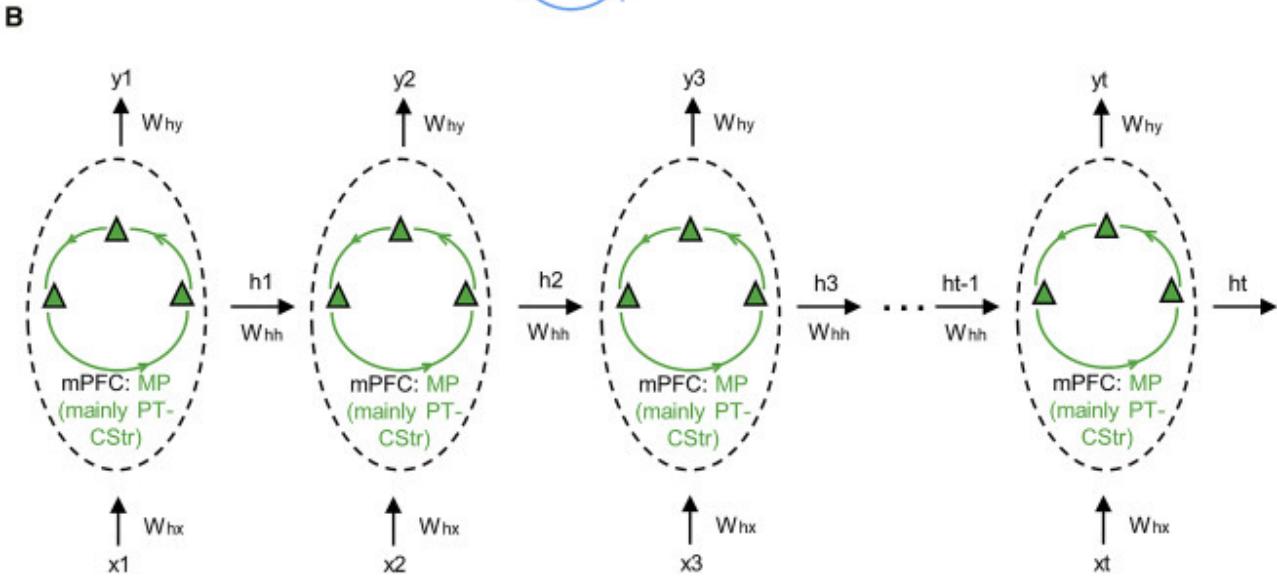
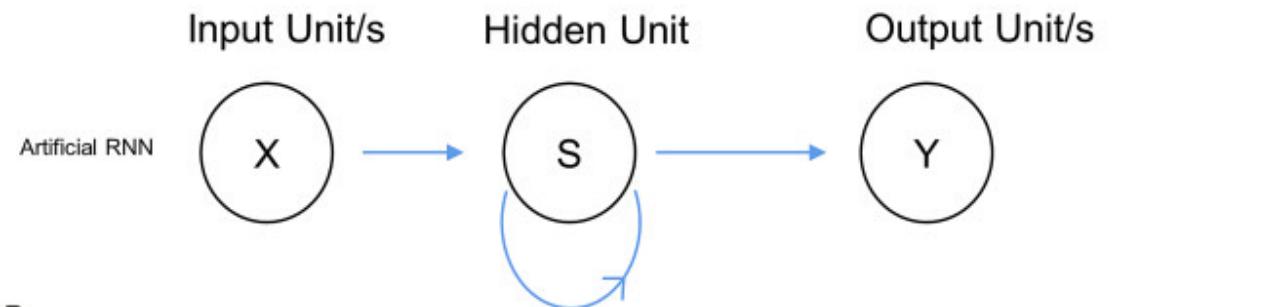
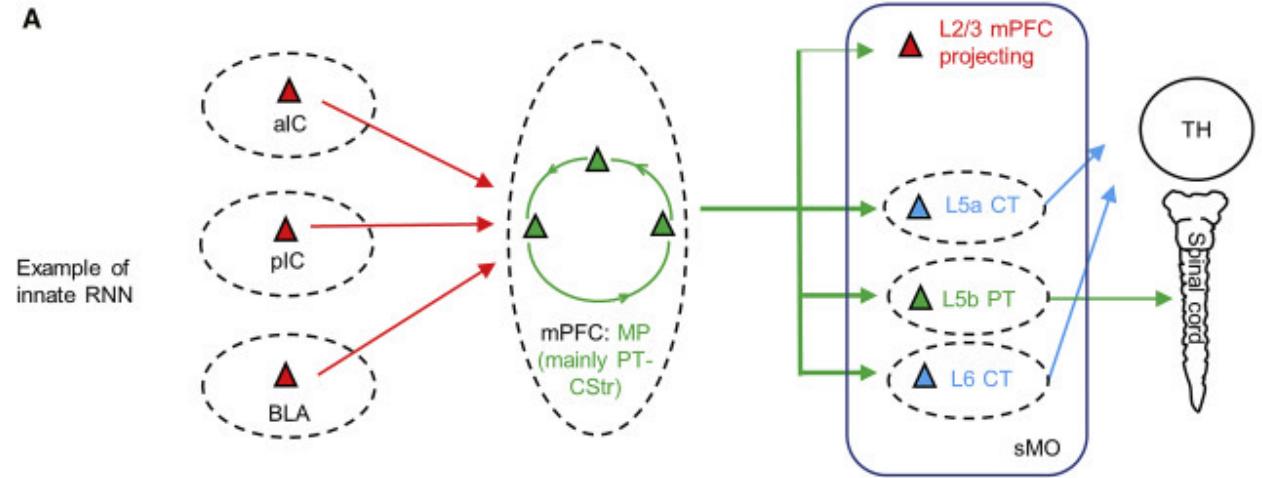
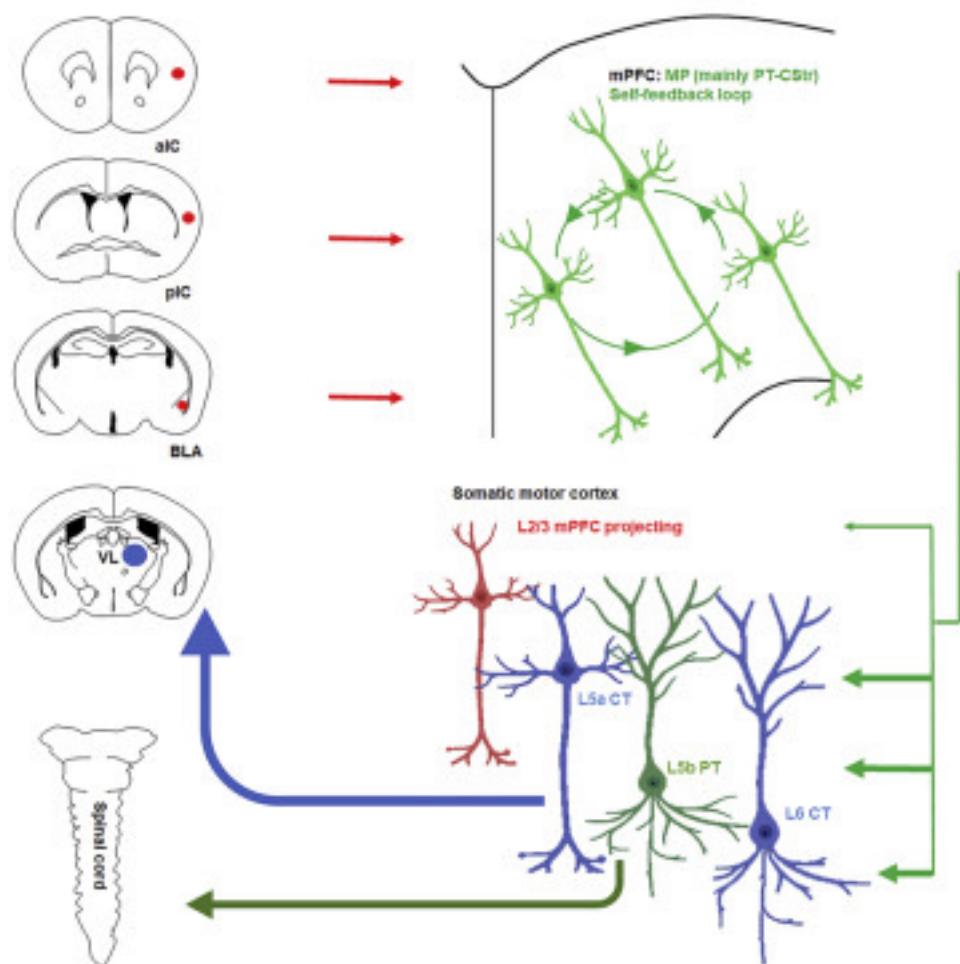


Figure 6

Basis Neurosains



Recurrent

$$\begin{aligned}\check{v}^{(1)} &= \check{v}^{(0)} + g\Delta t = g \\ \check{v}^{(2)} &= \check{v}^{(1)} + g\Delta t = g + g = 2g\end{aligned}$$

- Misalkan $s^{(t)}$ adalah situasi (state) sebuah sistem dinamik pada waktu t

- Dan θ adalah parameter lainnya pada fungsi $f()$

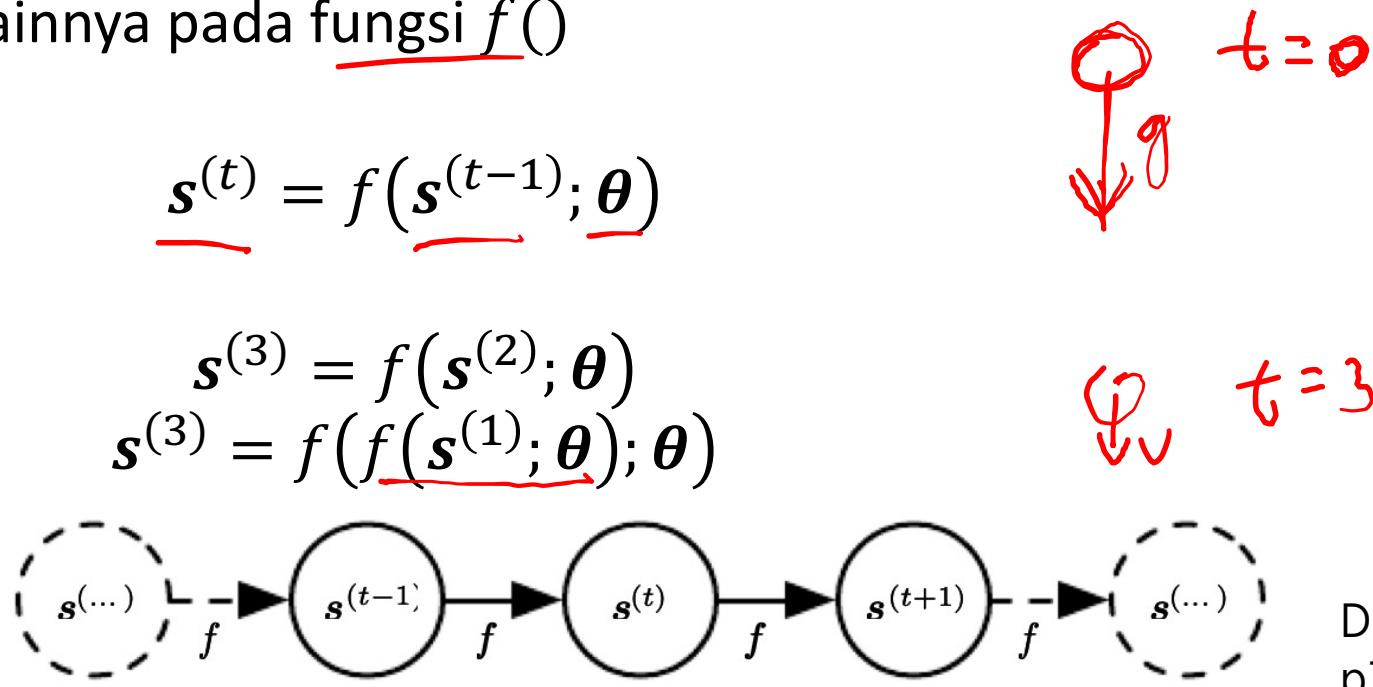
- Recurrent berarti:

$$\underline{s}^{(t)} = f(\underline{s}^{(t-1)}; \underline{\theta})$$

- Contoh:

$$\begin{aligned}s^{(3)} &= f(s^{(2)}; \theta) \\ s^{(3)} &= f(f(s^{(1)}; \theta); \theta)\end{aligned}$$

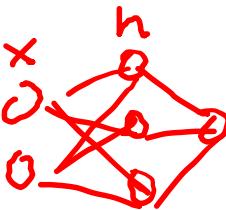
- Jejaring komputasi:



Deep Learning,
p375

Figure 10.1: The classical dynamical system described by equation 10.1, illustrated as an unfolded computational graph. Each node represents the state at some time t and the function f maps the state at t to the state at $t + 1$. The same parameters (the same value of θ used to parametrize f) are used for all time steps.

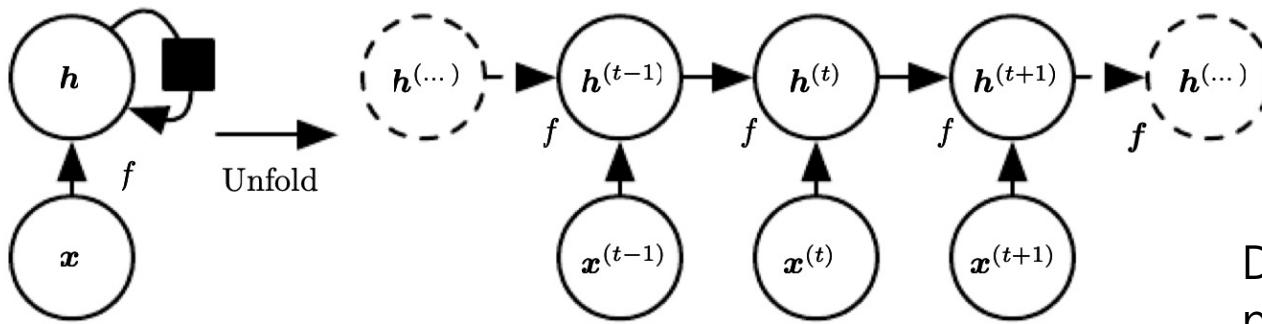
RNN



- Misalkan $\underline{h}^{(t)}$ adalah hidden unit sebuah neural network pada urutan t
- Dan $\underline{x}^{(t)}$ adalah external input terhadap neural network pada urutan t
- Dan $\underline{\theta}$ adalah parameter lainnya pada fungsi $f()$
- Recurrent Neural Network berarti:

$$\underline{h}^{(t)} = f(\underline{h}^{(t-1)}, \underline{x}^{(t)}; \underline{\theta})$$

- Jejaring komputasi:



Deep Learning,
p376

Figure 10.2: A recurrent network with no outputs. This recurrent network just processes information from the input x by incorporating it into the state h that is passed forward through time. (*Left*)Circuit diagram. The black square indicates a delay of a single time step. (*Right*)The same network seen as an unfolded computational graph, where each node is now associated with one particular time instance.

Kelebihan RNN

- Ukuran input selalu sama, tidak peduli kalimat sependek atau sepanjang apapun
- Menggunakan fungsi transisi dengan parameter yang sama

RNN Forward

- $\underline{a}^{(t)} = \underline{b} + \underline{W}h^{(t-1)} + \underline{U}x^{(t)}$
- $\underline{h}^{(t)} = \tanh(\underline{a}^{(t)})$
- $\underline{o}^{(t)} = \underline{c} + \underline{V}h^{(t)}$
- $\hat{\underline{y}}^{(t)} = \text{softmax}(\underline{o}^{(t)})$

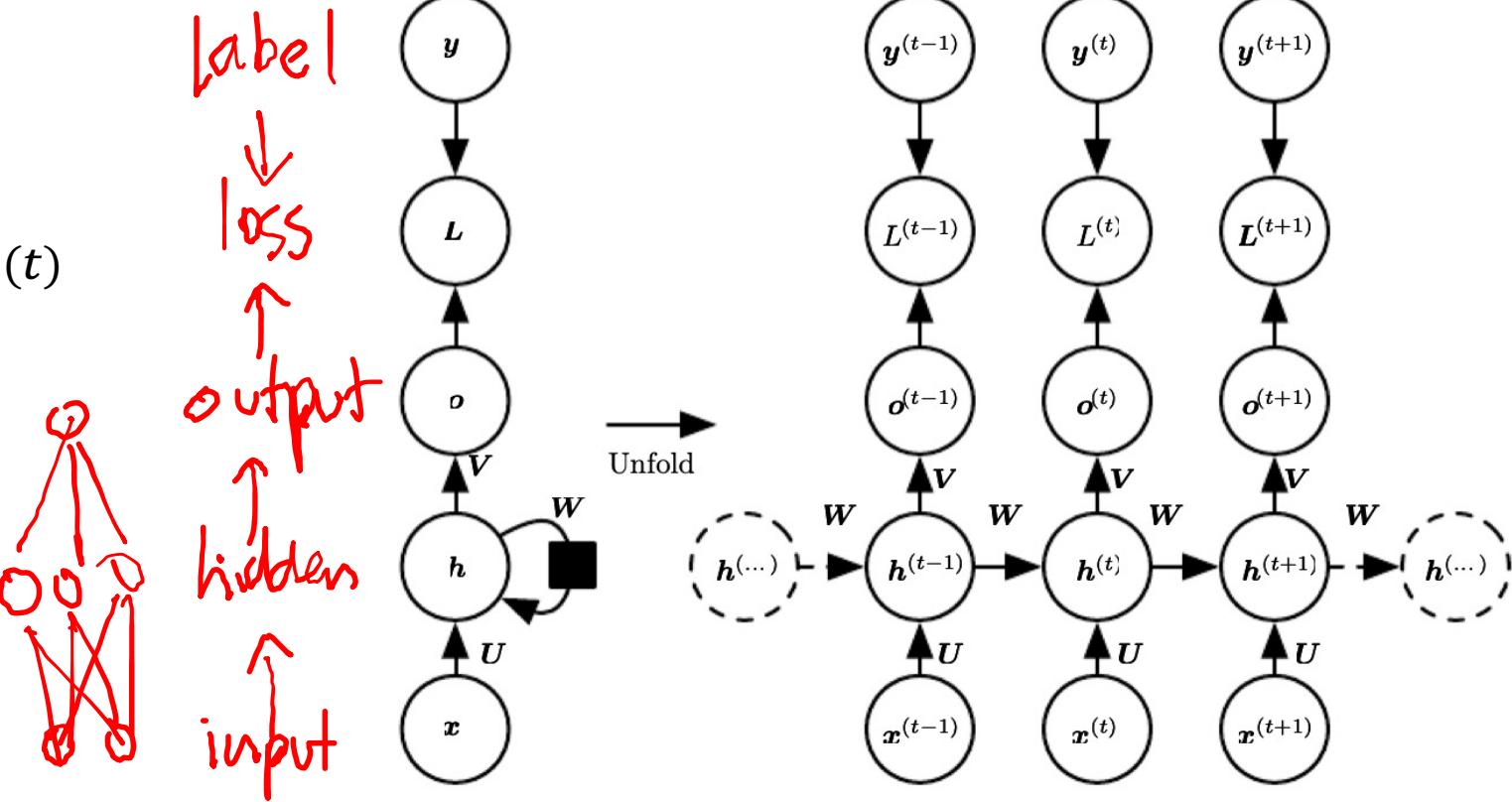
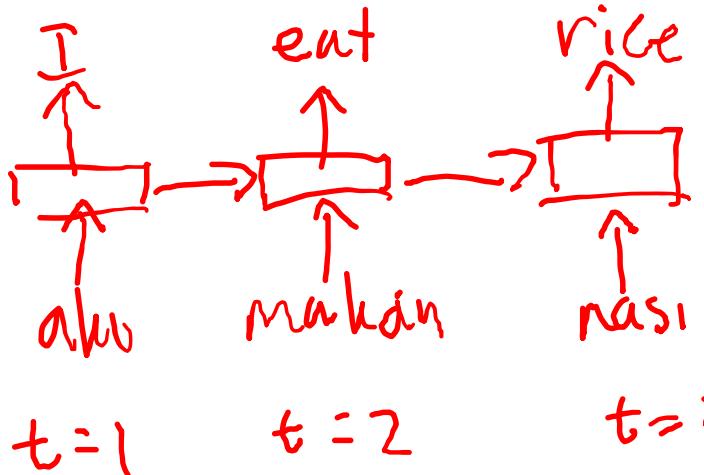
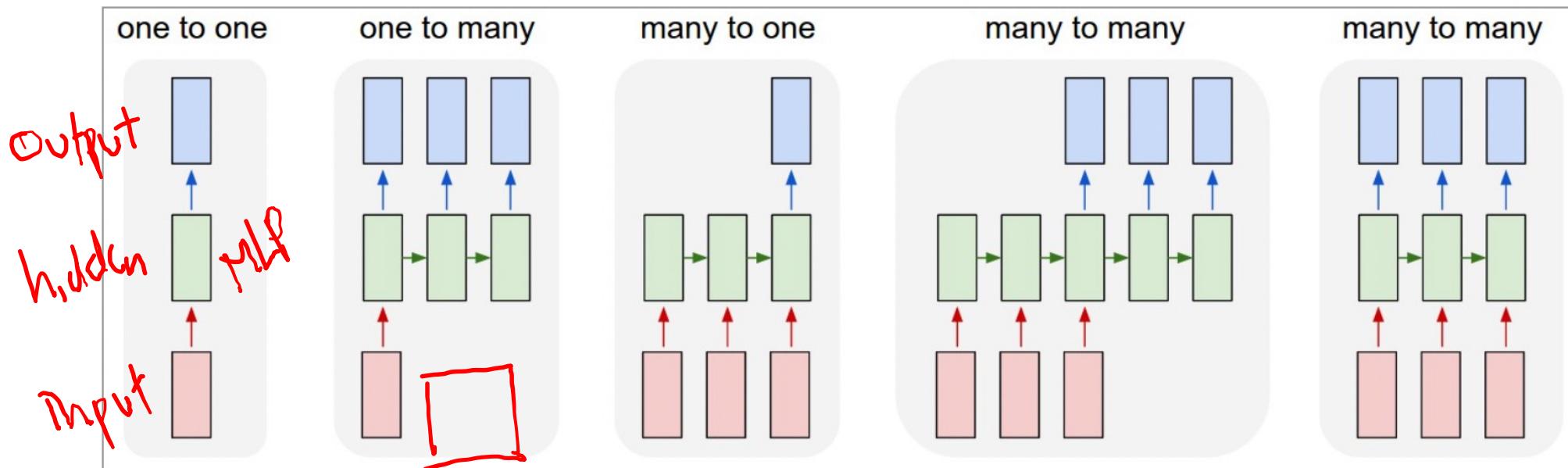


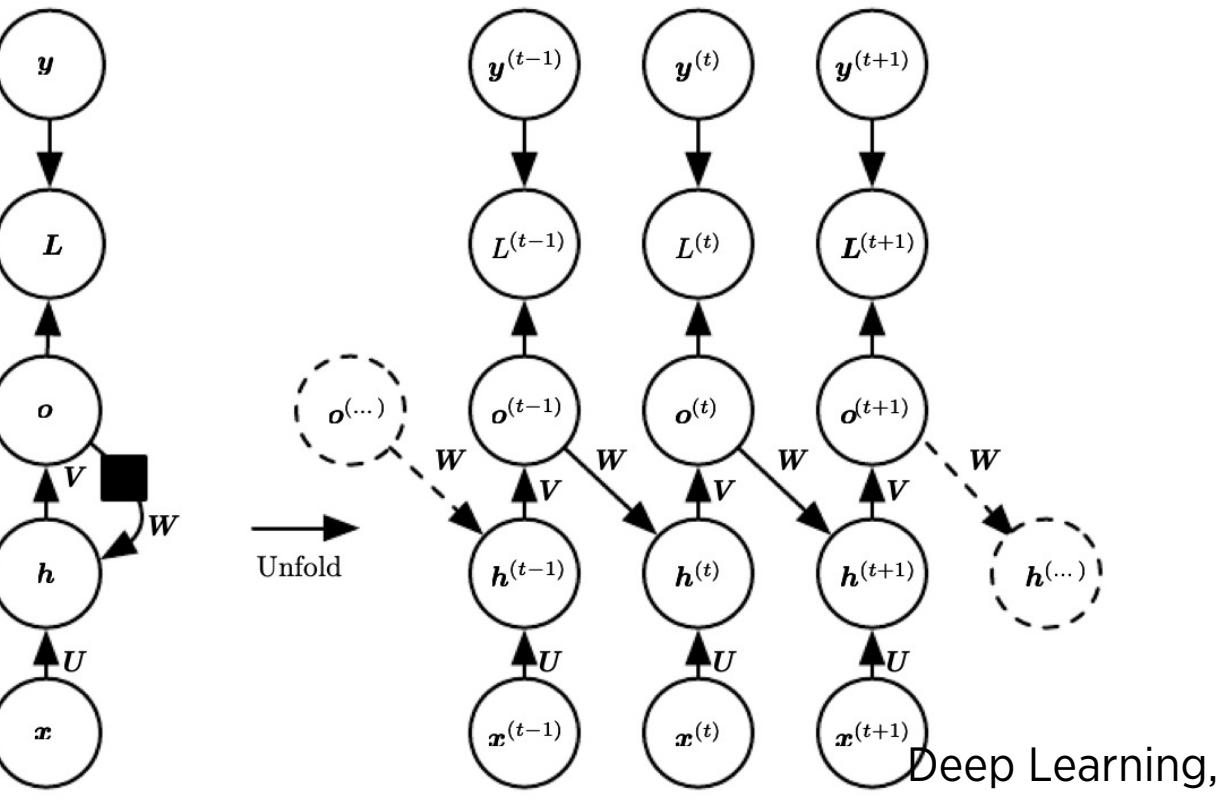
Figure 10.3: The computational graph to compute the training loss of a recurrent network that maps an input sequence of \mathbf{x} values to a corresponding sequence of output \mathbf{o} values. A loss L measures how far each \mathbf{o} is from the corresponding training target \mathbf{y} . When using softmax outputs, we assume \mathbf{o} is the unnormalized log probabilities. The loss L internally computes $\hat{\mathbf{y}} = \text{softmax}(\mathbf{o})$ and compares this to the target \mathbf{y} . The RNN has input to hidden connections parametrized by a weight matrix \mathbf{U} , hidden-to-hidden recurrent connections parametrized by a weight matrix \mathbf{W} , and hidden-to-output connections parametrized by a weight matrix \mathbf{V} . Equation 10.8 defines forward propagation in this model. (Left) The RNN and its loss drawn with recurrent connections. (Right) The same seen as an time-unfolded computational graph, where each node is now associated with one particular time instance.

Ragam input-output RNN



Each rectangle is a vector and arrows represent functions (e.g. matrix multiply). Input vectors are in red, output vectors are in blue and green vectors hold the RNN's state (more on this soon). From left to right: **(1)** Vanilla mode of processing without RNN, from fixed-sized input to fixed-sized output (e.g. image classification). **(2)** Sequence output (e.g. image captioning takes an image and outputs a sentence of words). **(3)** Sequence input (e.g. sentiment analysis where a given sentence is classified as expressing positive or negative sentiment). **(4)** Sequence input and sequence output (e.g. Machine Translation: an RNN reads a sentence in English and then outputs a sentence in French). **(5)** Synced sequence input and output (e.g. video classification where we wish to label each frame of the video). Notice that in every case are no pre-specified constraints on the lengths of sequences because the recurrent transformation (green) is fixed and can be applied as many times as we like.

Output to Input



Deep Learning,
p380

Figure 10.4: An RNN whose only recurrence is the feedback connection from the output to the hidden layer. At each time step t , the input is \mathbf{x}_t , the hidden layer activations are $\mathbf{h}^{(t)}$, the outputs are $\mathbf{o}^{(t)}$, the targets are $\mathbf{y}^{(t)}$ and the loss is $L^{(t)}$. (Left) Circuit diagram. (Right) Unfolded computational graph. Such an RNN is less powerful (can express a smaller set of functions) than those in the family represented by figure 10.3. The RNN in figure 10.3 can choose to put any information it wants about the past into its hidden representation \mathbf{h} and transmit \mathbf{h} to the future. The RNN in this figure is trained to put a specific output value into \mathbf{o} , and \mathbf{o} is the only information it is allowed to send to the future. There are no direct connections from \mathbf{h} going forward. The previous \mathbf{h} is connected to the present only indirectly, via the predictions it was used to produce. Unless \mathbf{o} is very high-dimensional and rich, it will usually lack important information from the past. This makes the RNN in this figure less powerful, but it may be easier to train because each time step can be trained in isolation from the others, allowing greater parallelization during training, as described in section 10.2.1.

Many to One RNN

- Biasa struktur seperti ini dapat digunakan untuk klasifikasi
- Misal: analisa sentimen, prediksi penulis

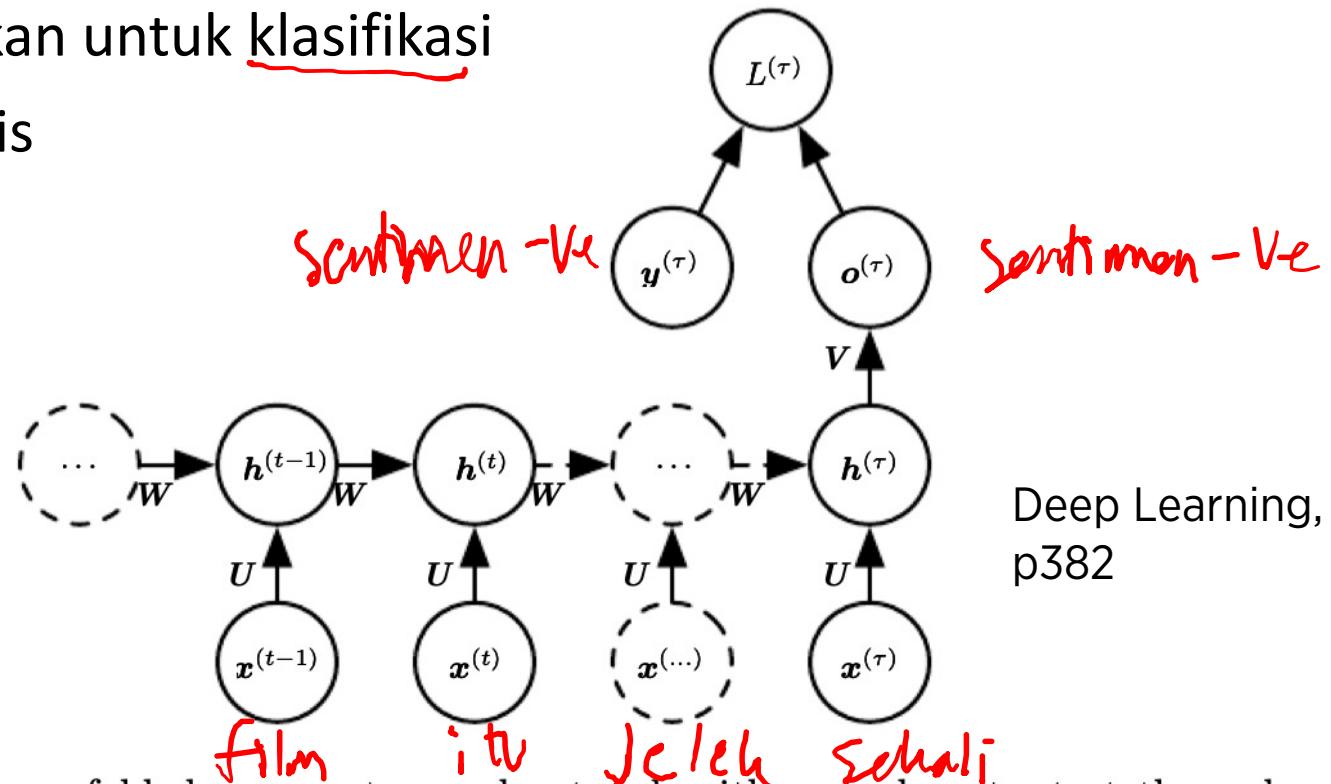
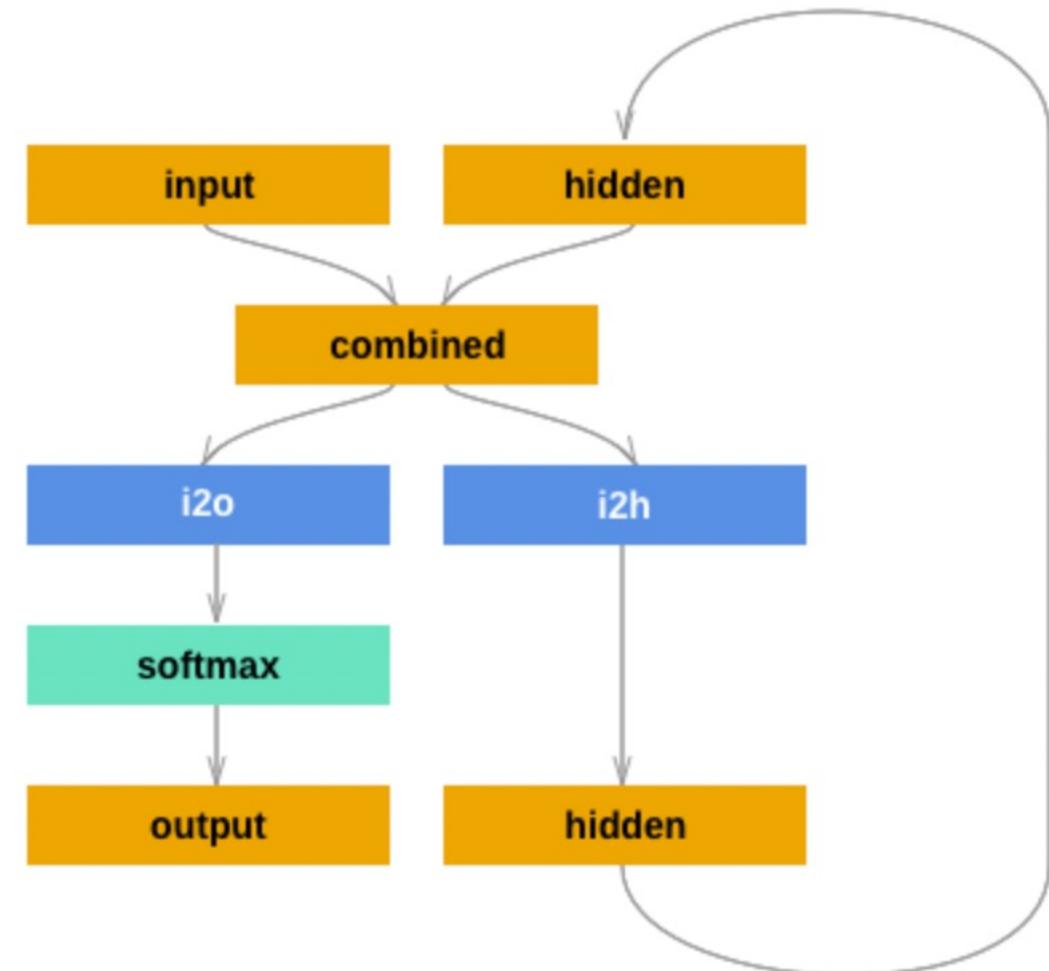


Figure 10.5: Time-unfolded recurrent neural network with a single output at the end of the sequence. Such a network can be used to summarize a sequence and produce a fixed-size representation used as input for further processing. There might be a target right at the end (as depicted here) or the gradient on the output $o^{(\tau)}$ can be obtained by back-propagating from further downstream modules.

Aktivitas: RNN di PyTorch

- Buatlah RNN dengan struktur berikut
- Gunakan RNN ini untuk prediksi bahasa



Menghitung Gradien

Using this notation, the gradient on the remaining parameters is given by:

$$\nabla_{\mathbf{c}} L = \sum_t \left(\frac{\partial \mathbf{o}^{(t)}}{\partial \mathbf{c}} \right)^\top \nabla_{\mathbf{o}^{(t)}} L = \sum_t \nabla_{\mathbf{o}^{(t)}} L \quad (10.22)$$

$$\nabla_{\mathbf{b}} L = \sum_t \left(\frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{b}^{(t)}} \right)^\top \nabla_{\mathbf{h}^{(t)}} L = \sum_t \text{diag} \left(1 - (\mathbf{h}^{(t)})^2 \right) \nabla_{\mathbf{h}^{(t)}} L \quad (10.23)$$

$$\nabla_{\mathbf{V}} L = \sum_t \sum_i \left(\frac{\partial L}{\partial o_i^{(t)}} \right) \nabla_{\mathbf{V}} o_i^{(t)} = \sum_t (\nabla_{\mathbf{o}^{(t)}} L) \mathbf{h}^{(t)\top} \quad (10.24)$$

$$\nabla_{\mathbf{W}} L = \sum_t \sum_i \left(\frac{\partial L}{\partial h_i^{(t)}} \right) \nabla_{\mathbf{W}^{(t)}} h_i^{(t)} \quad (10.25)$$

$$= \sum_t \text{diag} \left(1 - (\mathbf{h}^{(t)})^2 \right) (\nabla_{\mathbf{h}^{(t)}} L) \mathbf{h}^{(t-1)\top} \quad (10.26)$$

$$\nabla_{\mathbf{U}} L = \sum_t \sum_i \left(\frac{\partial L}{\partial h_i^{(t)}} \right) \nabla_{\mathbf{U}^{(t)}} h_i^{(t)} \quad (10.27)$$

$$= \sum_t \text{diag} \left(1 - (\mathbf{h}^{(t)})^2 \right) (\nabla_{\mathbf{h}^{(t)}} L) \mathbf{x}^{(t)\top} \quad (10.28)$$

We do not need to compute the gradient with respect to $\mathbf{x}^{(t)}$ for training because it does not have any parameters as ancestors in the computational graph defining the loss.

Bidirectional RNN

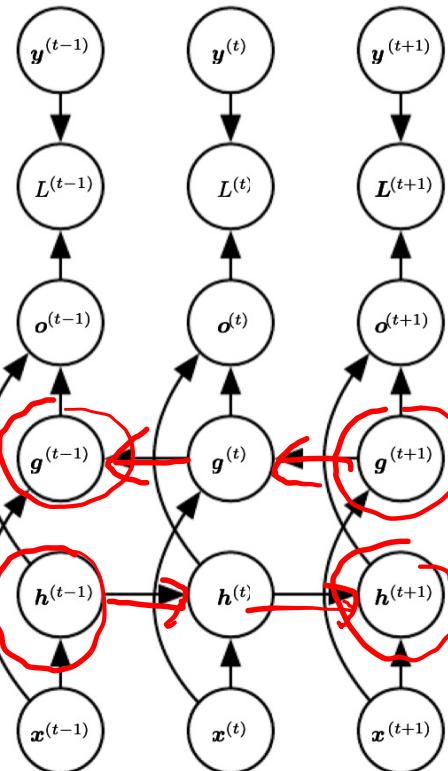


Figure 10.11: Computation of a typical bidirectional recurrent neural network, meant to learn to map input sequences \mathbf{x} to target sequences \mathbf{y} , with loss $L^{(t)}$ at each step t . The \mathbf{h} recurrence propagates information forward in time (towards the right) while the \mathbf{g} recurrence propagates information backward in time (towards the left). Thus at each point t , the output units $\mathbf{o}^{(t)}$ can benefit from a relevant summary of the past in its $\mathbf{h}^{(t)}$ input and from a relevant summary of the future in its $\mathbf{g}^{(t)}$ input.

Encoder-Decoder

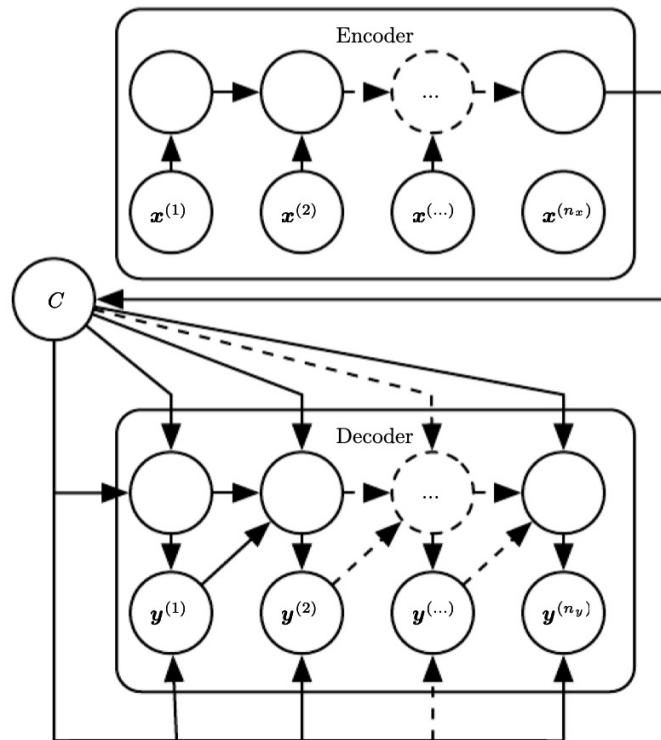
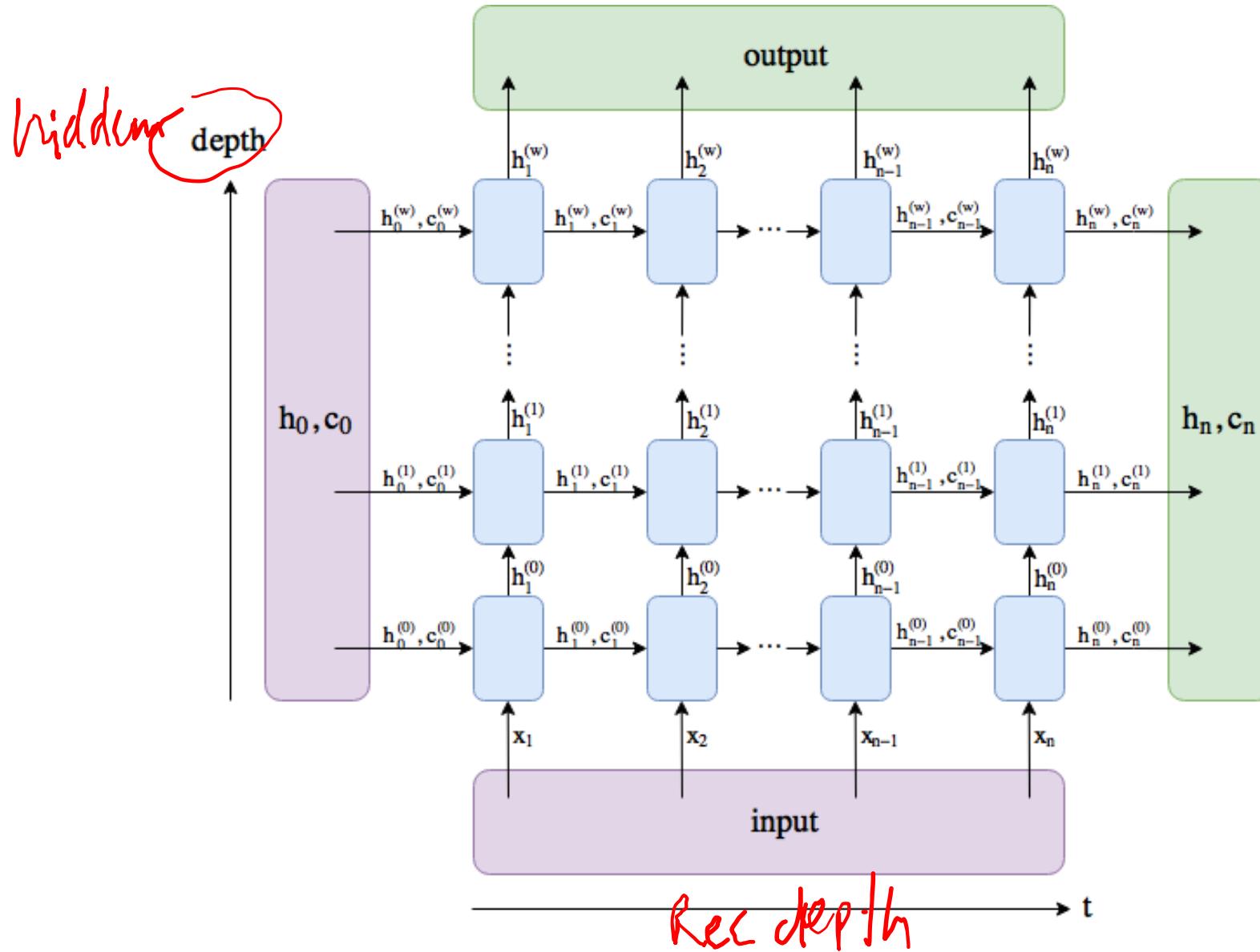


Figure 10.12: Example of an encoder-decoder or sequence-to-sequence RNN architecture, for learning to generate an output sequence $(y^{(1)}, \dots, y^{(n_y)})$ given an input sequence $(x^{(1)}, x^{(2)}, \dots, x^{(n_x)})$. It is composed of an encoder RNN that reads the input sequence and a decoder RNN that generates the output sequence (or computes the probability of a given output sequence). The final hidden state of the encoder RNN is used to compute a generally fixed-size context variable C which represents a semantic summary of the input sequence and is given as input to the decoder RNN.

Hidden Depth vs Recurrent Depth



Variasi Deep Recurrent Network

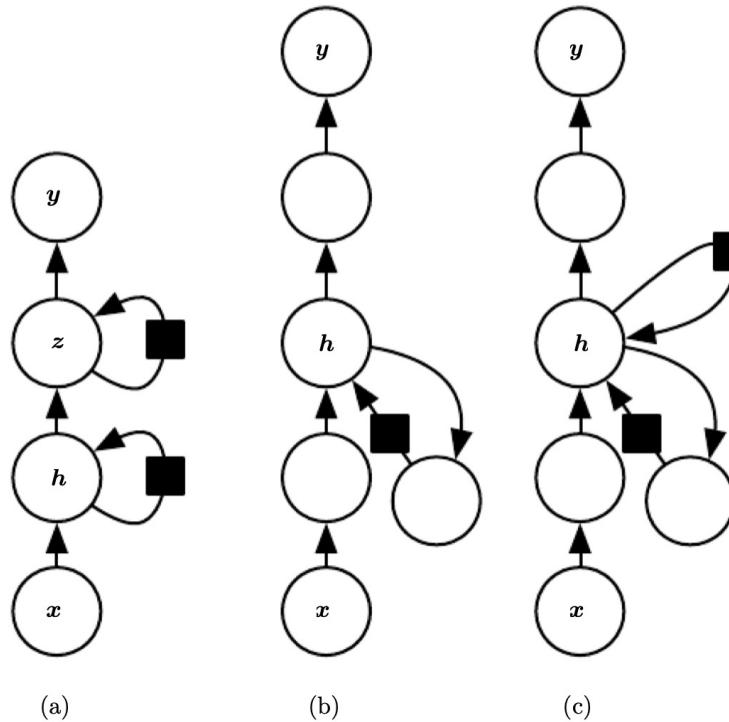


Figure 10.13: A recurrent neural network can be made deep in many ways (Pascanu *et al.*, 2014a). (a)The hidden recurrent state can be broken down into groups organized hierarchically. (b)Deeper computation (e.g., an MLP) can be introduced in the input-to-hidden, hidden-to-hidden and hidden-to-output parts. This may lengthen the shortest path linking different time steps. (c)The path-lengthening effect can be mitigated by introducing skip connections.

Recursive Network

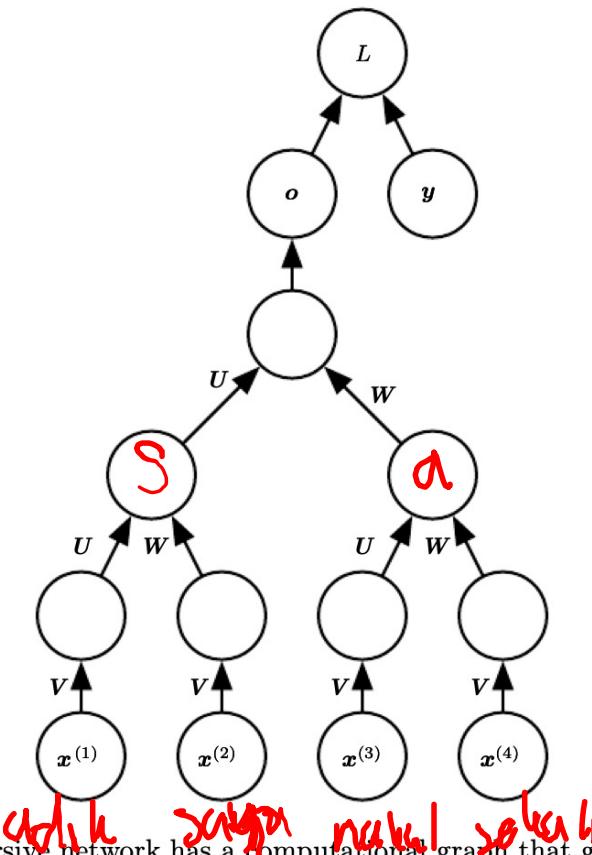


Figure 10.14: A recursive network has a computational graph that generalizes that of the recurrent network from a chain to a tree. A variable-size sequence $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(t)}$ can be mapped to a fixed-size representation (the output \mathbf{o}), with a fixed set of parameters (the weight matrices $\mathbf{U}, \mathbf{V}, \mathbf{W}$). The figure illustrates a supervised learning case in which some target \mathbf{y} is provided which is associated with the whole sequence.

LSTM & GRU

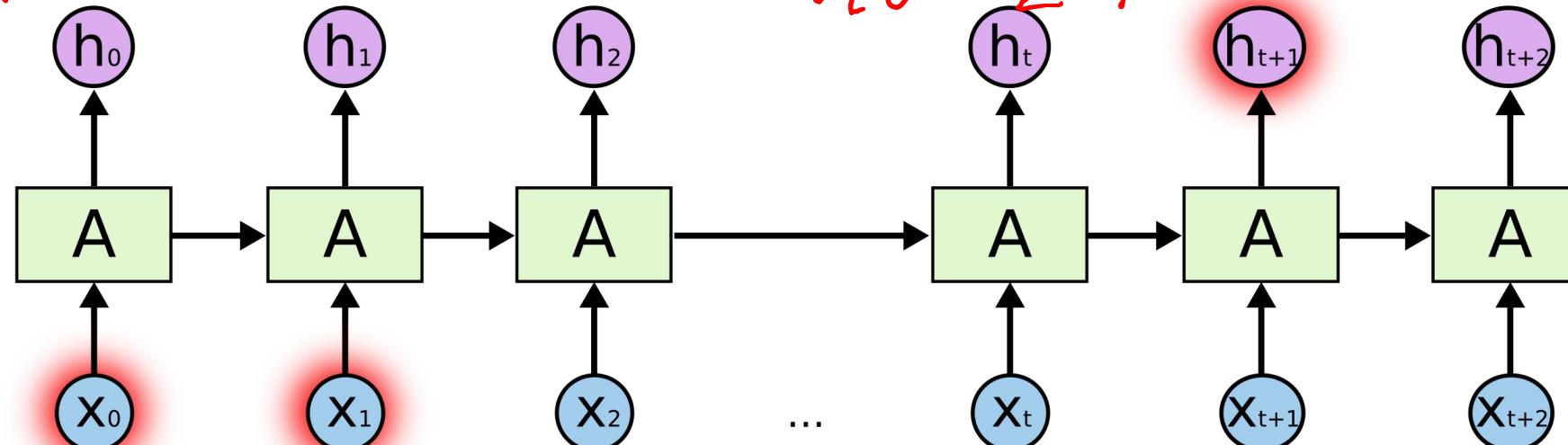
Vanishing/Exploding Gradient



- Hochreiter (1991) [German] and Bengio, et al. (1994) menemukan bahwa RNN tidak dapat mempelajari long term dependencies
- Contoh: "saya dibesarkan di Surabaya. Selama itu, saya ... Setelah itu, saya ... Saya bahkan masih lancar berbahasa Jawa."

$$W \times W \times W \times W \times \dots \times W$$

$$\begin{pmatrix} 0 & 0.1 \\ 0.1 & 0 \end{pmatrix} \begin{pmatrix} 0 & 0.1 \\ 0.1 & 0 \end{pmatrix} = \begin{pmatrix} 0 & 0.01 \\ 0.01 & 0 \end{pmatrix}$$



LSTM (Long Short-Term Memory)

$$0 < \sigma(x) < 1$$

Forget Gate

- $f_i^{(t)} = \sigma(b_i^f + \sum_j U_{i,j}^f x_j^{(t)} + \sum_j W_{i,j}^f h_j^{(t-1)})$

Input Gate

- $g_i^{(t)} = \sigma(b_i^g + \sum_j U_{i,j}^g x_j^{(t)} + \sum_j W_{i,j}^g h_j^{(t-1)})$

LSTM Cell

- $s_i^{(t)} = f_i^{(t)} s_i^{(t-1)} + g_i^{(t)} \sigma(b_i + \sum_j U_{i,j} x_j^{(t)} + \sum_j W_{i,j} h_j^{(t-1)})$

Output Gate

- $q_i^{(t)} = \sigma(b_i^o + \sum_j U_{i,j}^o x_j^{(t)} + \sum_j W_{i,j}^o h_j^{(t-1)})$
- $h_i^{(t)} = \tanh(s_i^{(t)}) q_i^{(t)}$

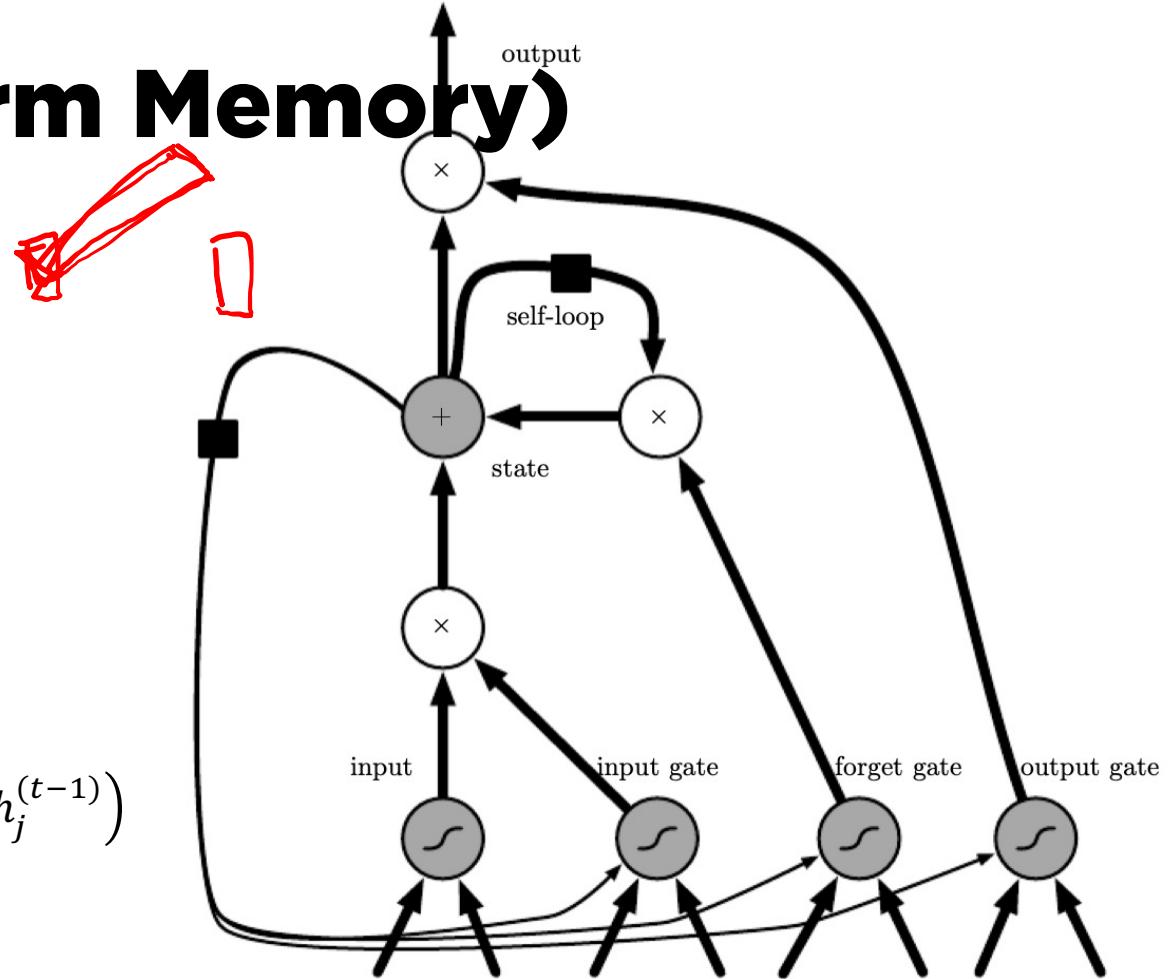


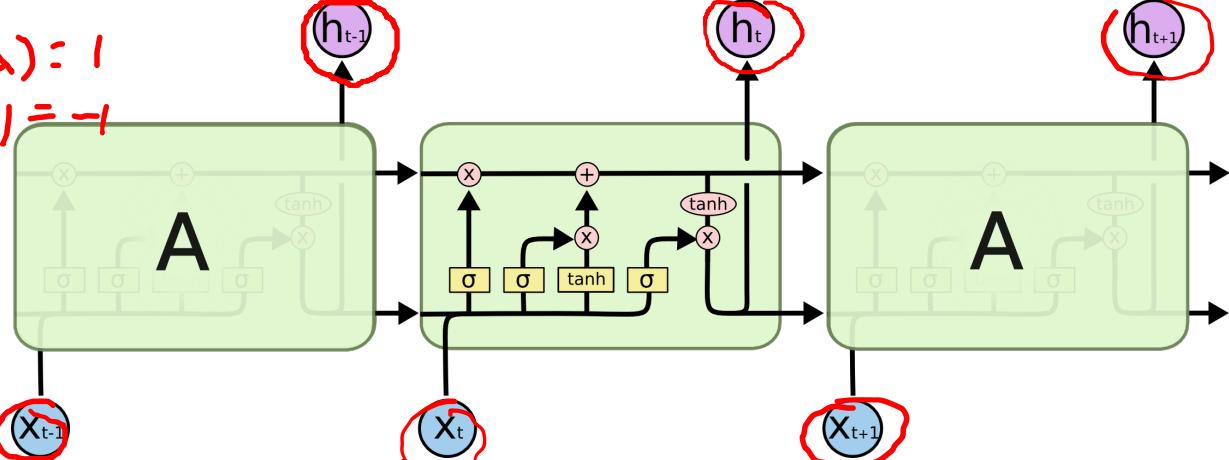
Figure 10.16: Block diagram of the LSTM recurrent network “cell.” Cells are connected recurrently to each other, replacing the usual hidden units of ordinary recurrent networks. An input feature is computed with a regular artificial neuron unit. Its value can be accumulated into the state if the sigmoidal input gate allows it. The state unit has a linear self-loop whose weight is controlled by the forget gate. The output of the cell can be shut off by the output gate. All the gating units have a sigmoid nonlinearity, while the input unit can have any squashing nonlinearity. The state unit can also be used as an extra input to the gating units. The black square indicates a delay of a single time step.

Skema

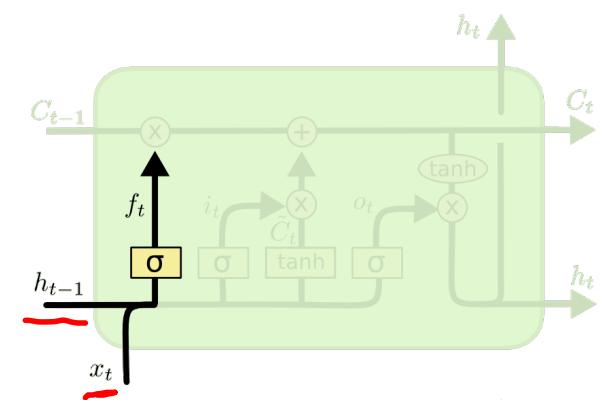
$$-1 < \tanh(x) < 1$$

$$x = 725 \rightarrow \tanh(x) = 1$$

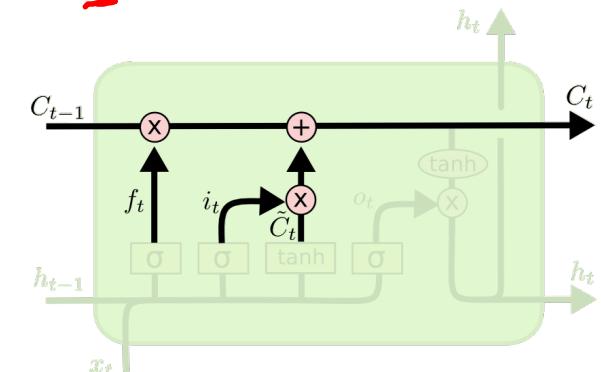
$$x = -953 \rightarrow \tanh(x) = -1$$



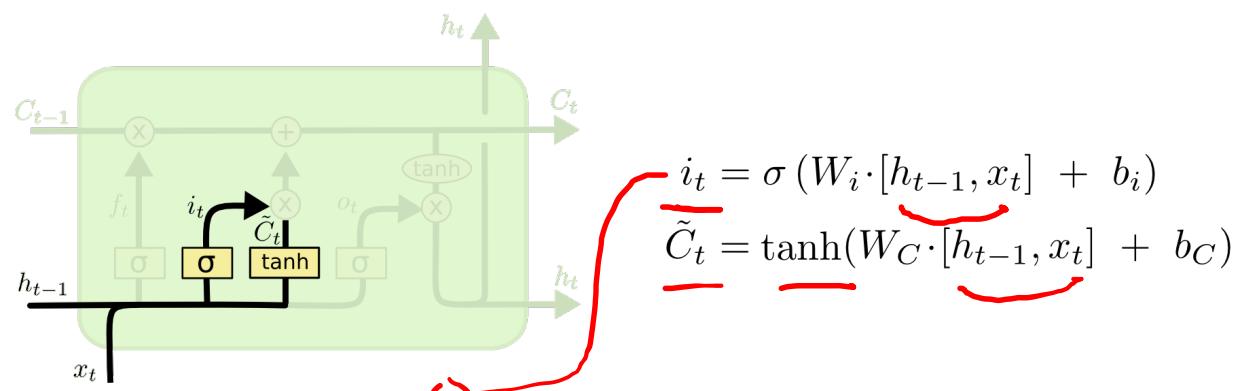
<https://colah.github.io/posts/2015-08-Understanding-LSTMs/>



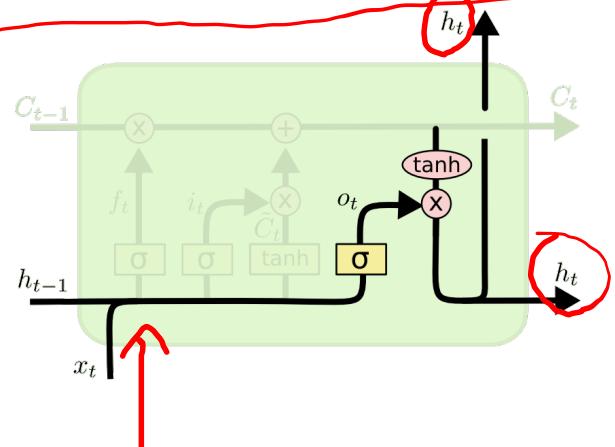
$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$



$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$



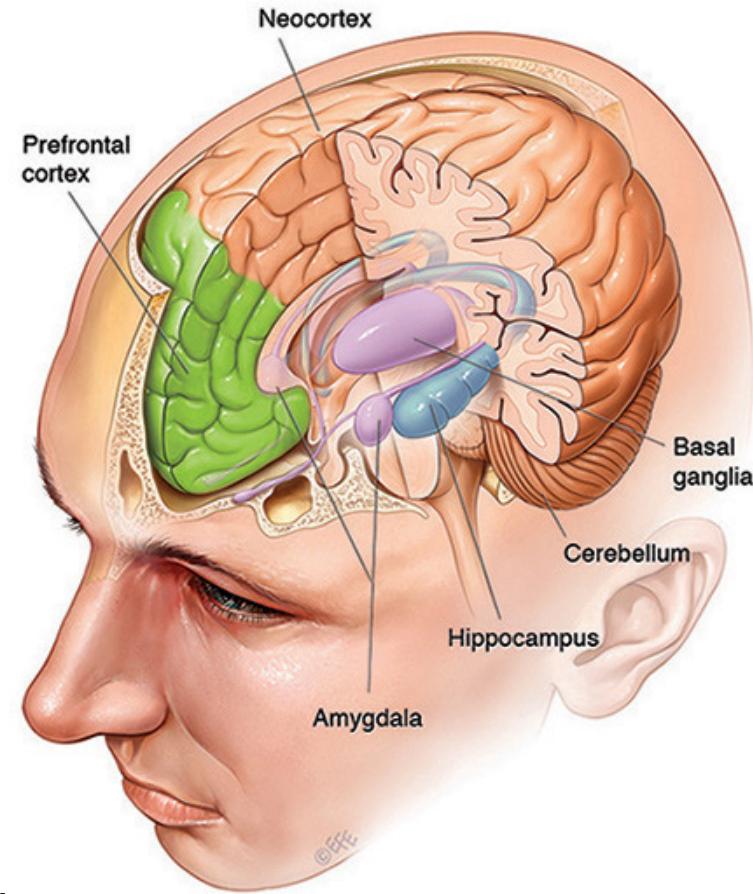
$$\begin{aligned} i_t &= \sigma(W_i \cdot [h_{t-1}, x_t] + b_i) \\ \tilde{C}_t &= \tanh(W_C \cdot [h_{t-1}, x_t] + b_C) \end{aligned}$$



$$\begin{aligned} o_t &= \sigma(W_o \cdot [h_{t-1}, x_t] + b_o) \\ h_t &= o_t * \tanh(C_t) \end{aligned}$$

Hippocampus

- In 1953, Henry Molaison had his hippocampus surgically removed to treat his epilepsy. His epilepsy was cured. However, after the surgery he was only able to form episodic memories that lasted a matter of minutes (completely unable to permanently store new information). As a result, Molaison's memory became mostly limited to events that occurred years before his surgery, in the distant past. He was, however, still able to improve his performance on various motor tasks, even though he had no memory of ever encountering or practising them. This indicated that although the hippocampus is crucial for laying down memories, it is not the site of permanent memory storage and isn't needed for motor memories (basal ganglia and cerebellum).

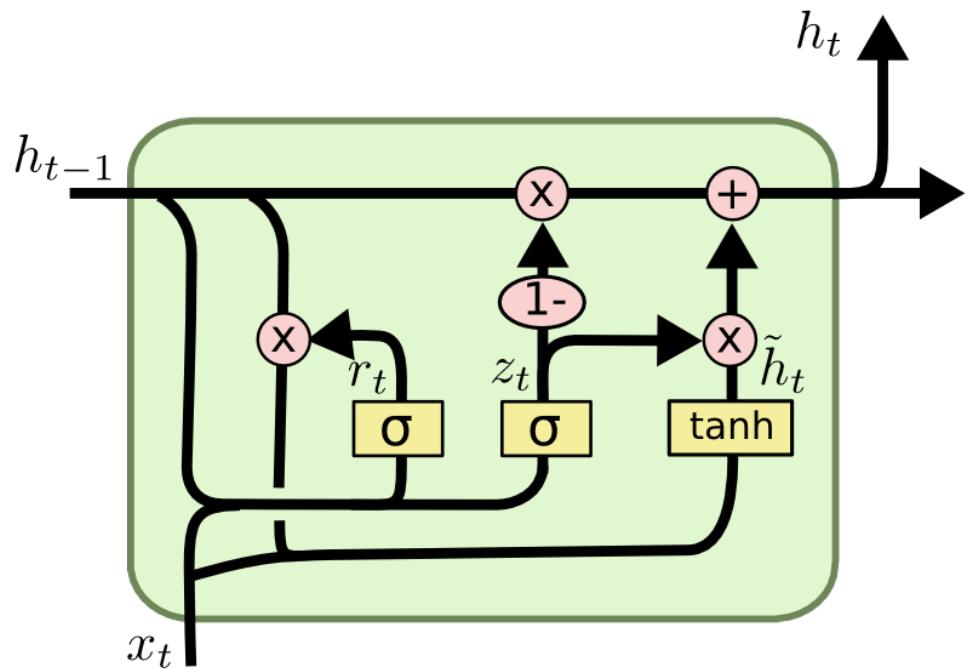


<https://www.telegraph.co.uk/culture/10047050/Henry-Molaison-The-incredible-story-of-the-man-with-no-memory.html>

GRU (Gated Recurrent Units)

- *Update Gate*
- $u_i^{(t)} = \sigma(b_i^u + \sum_j U_{i,j}^u x_j^{(t)} + \sum_j W_{i,j}^u h_j^{(t)})$
- *Reset Gate*
- $r_i^{(t)} = \sigma(b_i^r + \sum_j U_{i,j}^r x_j^{(t)} + \sum_j W_{i,j}^r h_j^{(t)})$
- *Update Equation*
- $h_i^{(t)} = u_i^{(t-1)} h_i^{(t-1)} + (1 - u_i^{(t-1)}) \sigma(b_i + \sum_j U_{i,j} x_j^{(t-1)} + \sum_j W_{i,j} r_j^{(t-1)} h_j^{(t-1)})$

Skema



$$z_t = \sigma(W_z \cdot [h_{t-1}, x_t])$$

$$r_t = \sigma(W_r \cdot [h_{t-1}, x_t])$$

$$\tilde{h}_t = \tanh(W \cdot [r_t * h_{t-1}, x_t])$$

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

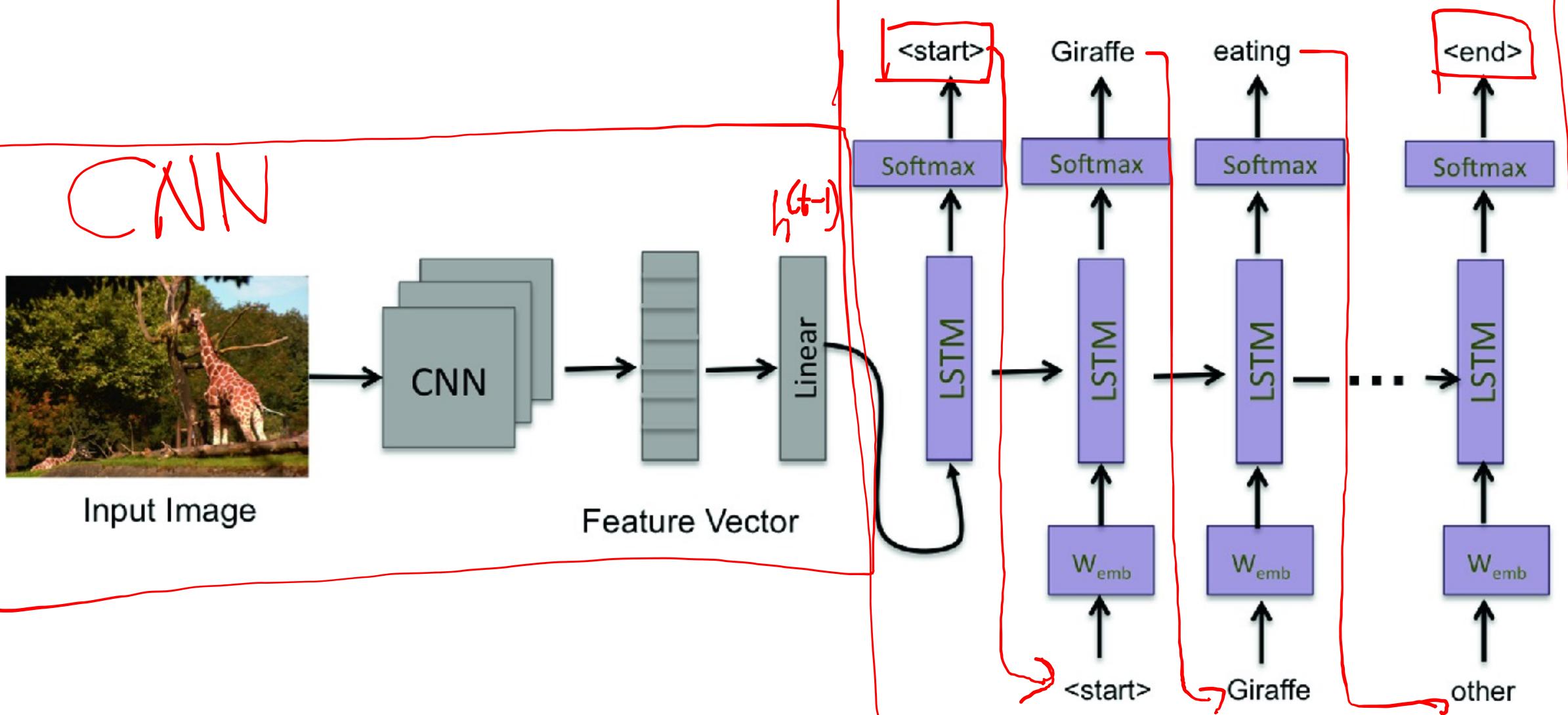
[]

[]

Aktivitas: Time-Series Menggunakan Sequence Model

- Gunakan sequence model untuk memprediksi pola time-series pada 2 dataset yang sudah disediakan:
 - Temperature Melbourne
 - Sunspots
- Gunaan 80% awal series sebagai data latih dan 20% akhir series sebagai data tes
- Bandingkan RNN, LSTM, GRU:
 - MAE
 - Latensi train dan test

Image Captioning



Summary

- Masing-masing sebutkan konsep utama yang diingat

Tuhan Memberkati

