

COMP4434 Big Data Analytics

Lecture 7

Recurrent Neural Networks

HUANG Xiao

xiaohuang@comp.polyu.edu.hk



Outline

- Vanilla Recurrent Neural Networks
 - Exploding and Vanishing Gradients
- RNNs with Cell states
 - Long Short-Term Memory (LSTM)
 - Gated Recurrent Unit (GRU)
- Sequence Learning Architectures
 - Sequence Learning with one RNN Layer
 - Sequence Learning with multiple RNN Layers
- Application: Sequence-to-Sequence Model in Activities of Daily Living (ADL) Recognition

Recurrent Neural Networks

- Human brain deals with information streams. Most data is obtained, processed, and generated sequentially.
 - E.g., listening: soundwaves → vocabularies/sentences
 - E.g., reading: words in sequence
- Human thoughts have persistence; humans don't start their thinking from scratch every second.
 - As you read this sentence, you understand each word based on your prior knowledge
- The applications of multilayer perceptron and Convolutional Neural Networks are limited due to:
 - Only accept a fixed-size vector as input (e.g., an image) and produce a fixed-size vector as output (e.g., probabilities of different classes)
- Recurrent Neural Networks (RNNs) are a family of neural networks introduced to **learn sequential data**.
 - Inspired by the temporal-dependent and persistent human thoughts

Real-life Sequence Learning Applications

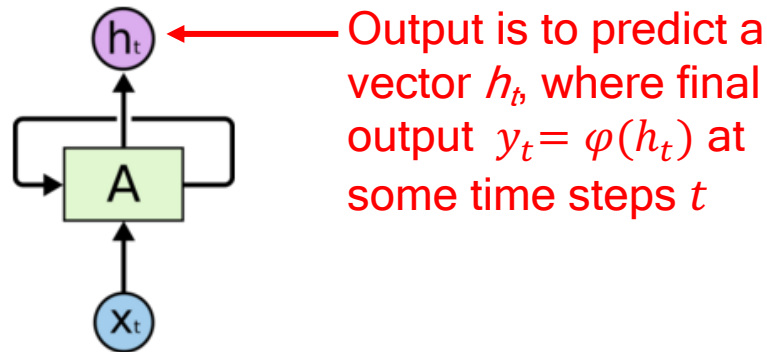
RNNs can be applied to various type of sequential data to learn the temporal patterns.

- Time-series data (e.g., stock price) → Prediction, regression
- Raw sensor data (e.g., signal, voice, handwriting) → Labels or text sequences
- Text → Label (e.g., sentiment) or text sequence (e.g., translation, summary, answer)
- Image and video → Text description (e.g., captions, scene interpretation)

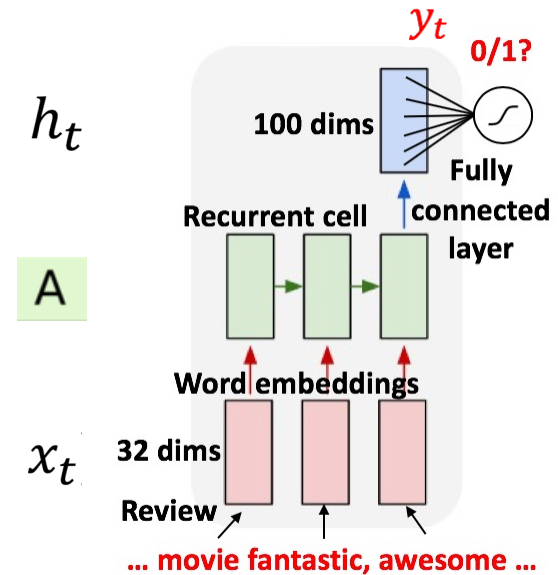
Task	Input	Output
Activity Recognition (Zhu et al. 2018)	Sensor Signals	Activity Labels
Machine translation (Sutskever et al. 2014)	English text	French text
Question answering (Bordes et al. 2014)	Question	Answer
Speech recognition (Graves et al. 2013)	Voice	Text
Handwriting prediction (Graves 2013)	Handwriting	Text
Opinion mining (Irsoy et al. 2014)	Text	Opinion expression

RNNs have loops

RNNs are networks with loops, allowing information to persist.



A chunk of neural network, $A = f_W$, looks at some input x_t and outputs a value h_t . A loop allows information to be passed from one step of the network to the next.



$$\boxed{h_t} = f_W(\boxed{h_{t-1}}, \boxed{x_t})$$

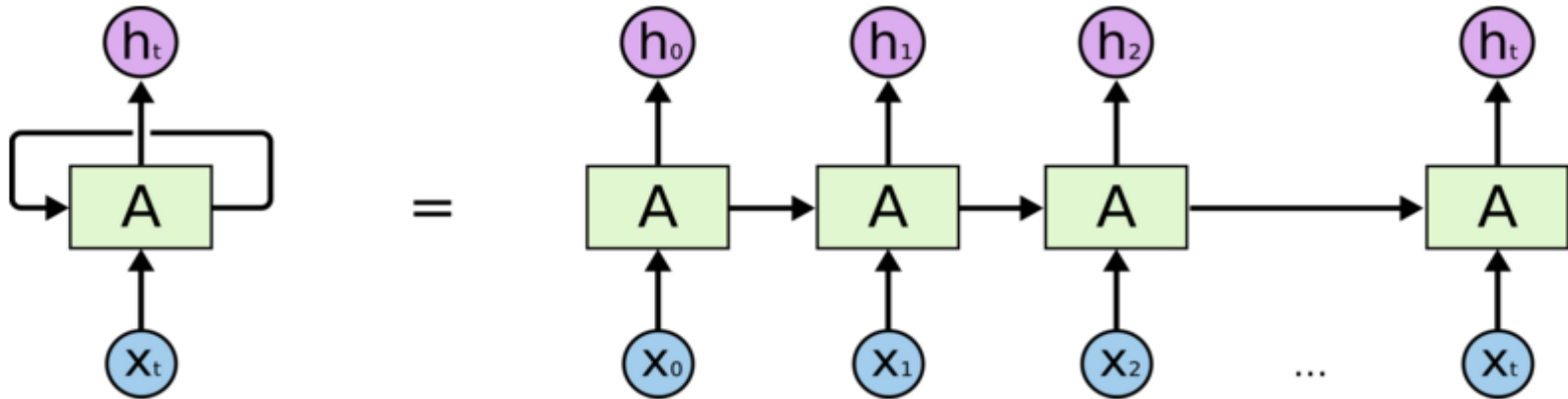
new state

old state

function with parameter W

Input vector at some time step

Unrolling RNN



An unrolled recurrent neural network.

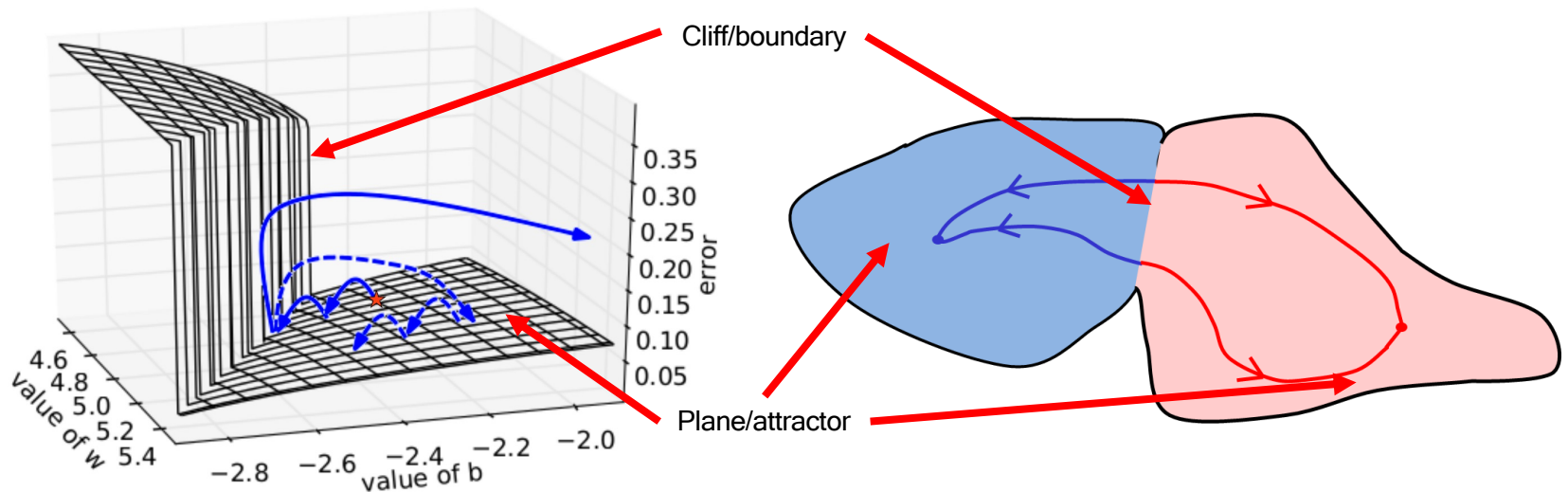
- A recurrent neural network can be thought of as multiple copies of the same network, each passing a message to a successor
- The diagram above shows what happens if we unroll the loop

Pros and cons of vanilla RNNs

The recurrent structure of RNNs enables the following characteristics:

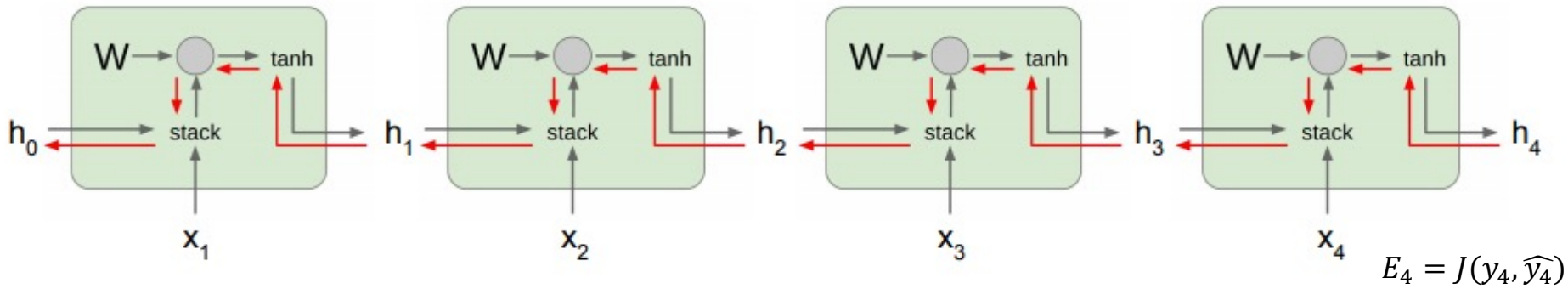
- Specialized for processing a sequence of values $x^{(1)}, \dots, x^{(\tau)}$
 - Each value $x^{(i)}$ is processed with the **same network \mathbf{A} that preserves past information**
- Can scale to much **longer sequences** than would be practical for networks without a recurrent structure
 - Reusing network **\mathbf{A}** reduces the required amount of parameters in the network
- Can process **variable-length sequences**
 - The network complexity does not vary when the input length change
- However, vanilla RNNs suffer from the training difficulty due to **exploding and vanishing gradients**

Exploding and Vanishing Gradients



- **Exploding:** If we start almost exactly on the boundary (cliff), tiny changes can make a huge difference.
- **Vanishing:** If we start a trajectory within an attractor (plane, flat surface), small changes in where we start make no difference to where we end up.
- Both cases hinder the learning process.

Exploding and Vanishing Gradients in RNNs



In vanilla RNNs, computing this gradient involves many factors of W (and repeated \tanh)*. If we decompose the singular values of the gradient multiplication matrix,

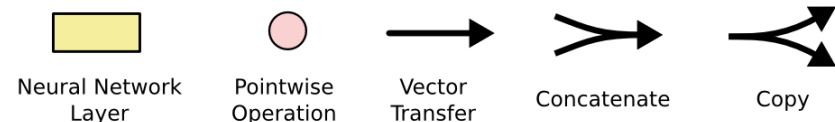
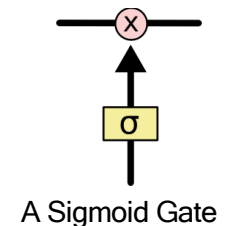
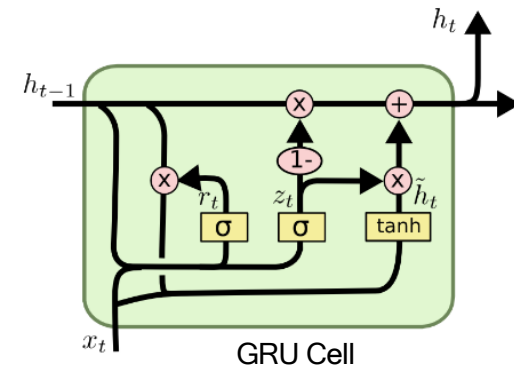
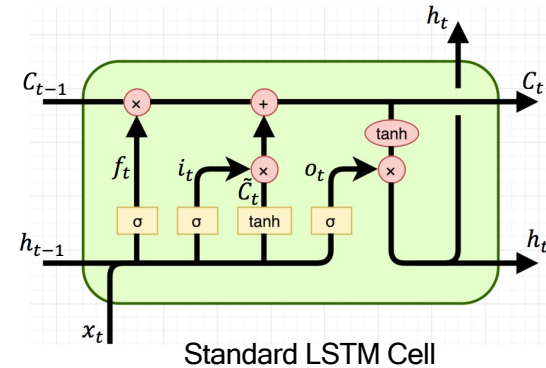
- Largest singular value $> 1 \rightarrow$ **Exploding gradients**
 - Slight error in the late time steps causes drastic updates in the early time steps \rightarrow Unstable learning
- Largest singular value $< 1 \rightarrow$ **Vanishing gradients**
 - Gradients passed to the early time steps is close to 0. \rightarrow Uninformed correction

* Refer to Bengio et al. (1994) or Goodfellow et al. (2016) for a complete derivation

RNNs with Cell states

- Vanilla RNN operates in a “multiplicative” way (repeated tanh)
- Two recurrent cell designs were proposed and widely adopted:
 - Long Short-Term Memory (LSTM)** (Hochreiter and Schmidhuber, 1997)
 - Gated Recurrent Unit (GRU) (Cho et al. 2014)
- Both designs process information in an “additive” way with gates to control information flow
 - Sigmoid gate outputs numbers between 0 and 1, describing how much of each component should be let through**

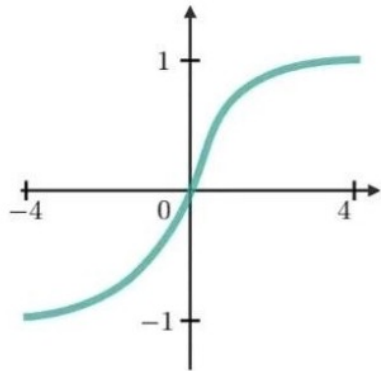
E.g., $f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f) = \text{Sigmoid}(W_f x_t + U_t h_{t-1} + b_f)$



Activation Functions

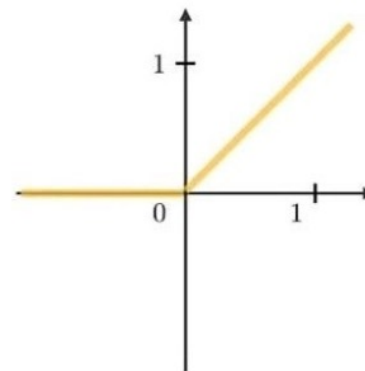
Tanh

$$g(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$



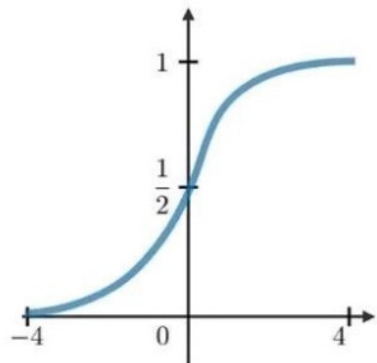
RELU

$$g(z) = \max(0, z)$$



Sigmoid

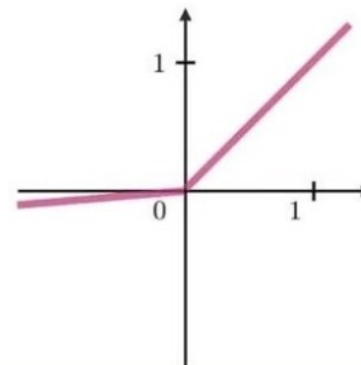
$$g(z) = \frac{1}{1 + e^{-z}}$$



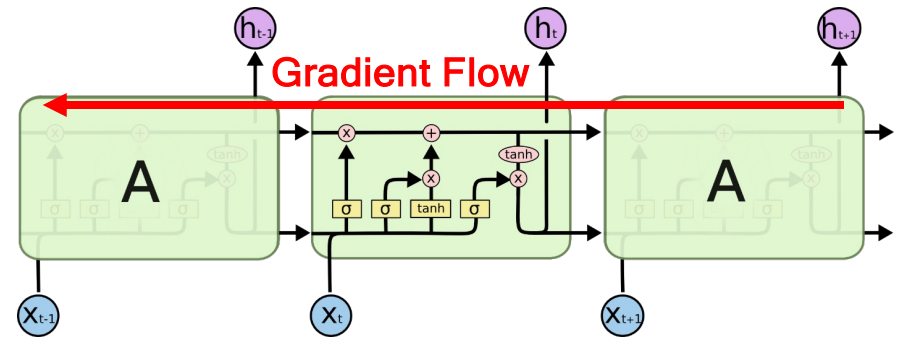
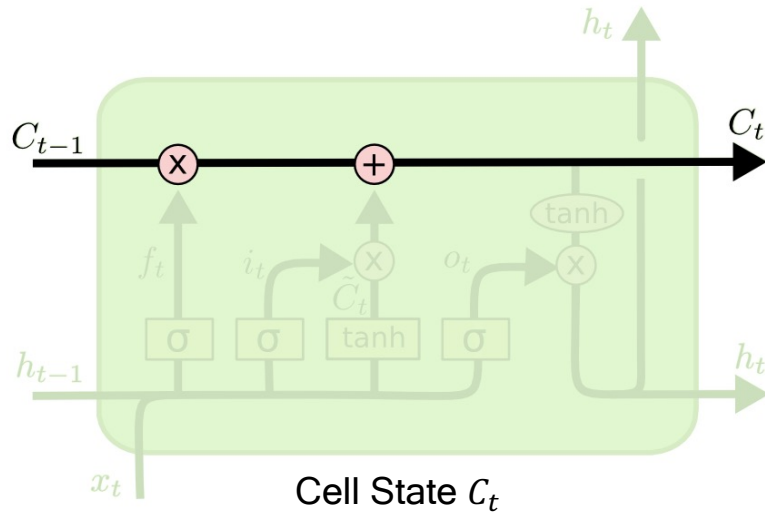
Leaky ReLU

$$g(z) = \max(\epsilon z, z)$$

with $\epsilon \ll 1$

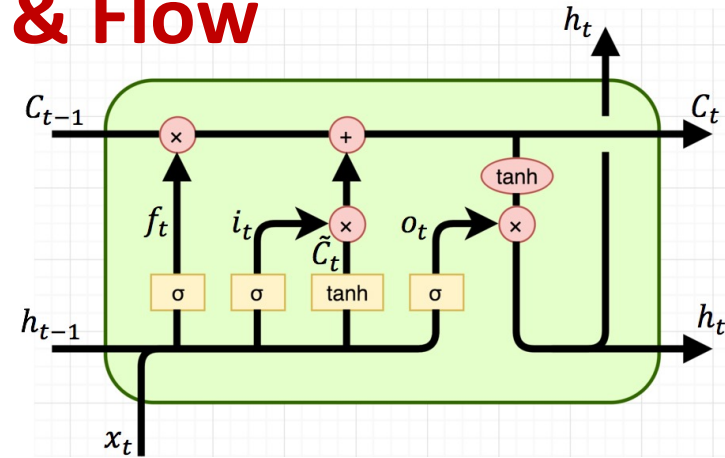


Long Short-Term Memory (LSTM)



- The key to LSTMs is the **cell state**.
 - Stores information of the past → long-term memory
 - Passes along time steps with minor linear interactions → “additive”
 - Results in an **uninterrupted gradient flow** → errors in the past pertain and impact learning in the future
- The LSTM cell manipulates input information with three gates.
 - **Input gate** → controls the intake of new information
 - **Forget gate** → determines what part of the cell state to be updated
 - **Output gate** → determines what part of the cell state to output

LSTM: Components & Flow



- LSTM unit output
- **Output gate** units
- Transformed memory cell contents
- Gated update to memory cell units
- **Forget gate** units
- **Input gate** units
- Potential input to memory cell

$$h_t = o_t * \tanh(C_t)$$

$$o_t = \sigma(W_o [h_{t-1}, x_t] + b_o)$$

$$\tanh(C_t)$$

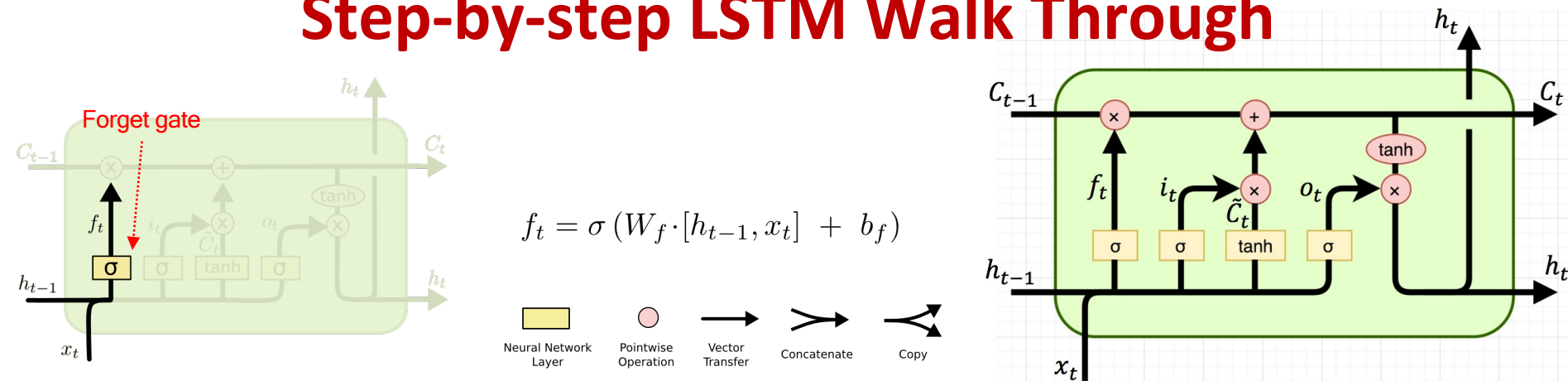
$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

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

$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

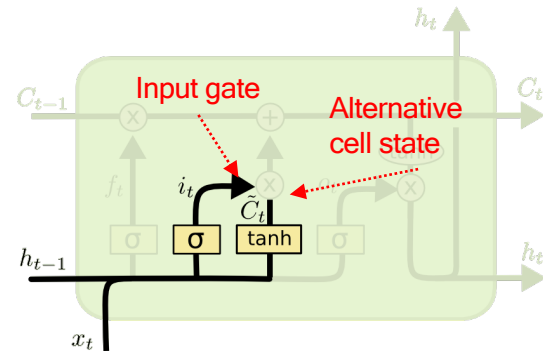
Step-by-step LSTM Walk Through



Step 1: Decide what information to throw away from the cell state (memory) $\Rightarrow f_t * C_{t-1}$.

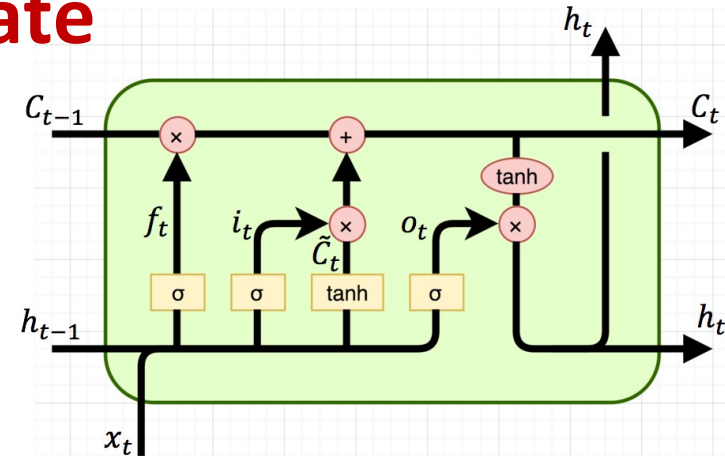
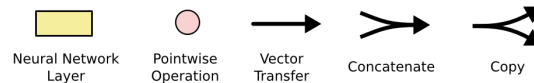
- The output of the previous state h_{t-1} and the new information x_t jointly determine what to forget
 - h_{t-1} contains selected features from the memory C_{t-1}
- Forget gate f_t ranges between $[0, 1]$
- Text processing example:
 - Cell state may include the gender of the current subject (h_{t-1}).
 - When the model observes a new subject (x_t), it may want to forget ($f_t \rightarrow 0$) the old subject in the memory (C_{t-1}).

Step 2: Input gate



$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

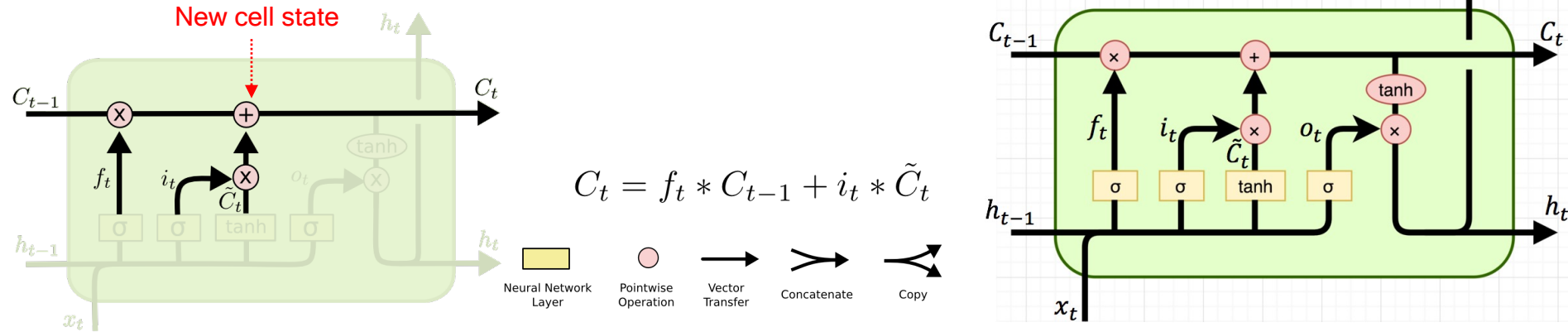
$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$



Step 2: Prepare the updates for the cell state from input $\rightarrow i_t * \tilde{C}_t$

- An alternative cell state \tilde{C}_t is created from the new information x_t with the guidance of h_{t-1}
- Input gate i_t ranges between $[0, 1]$
- **Example:** the model may want to add ($i_t \rightarrow 1$) the gender of new subject (\tilde{C}_t) to the cell state to replace the old one it is forgetting

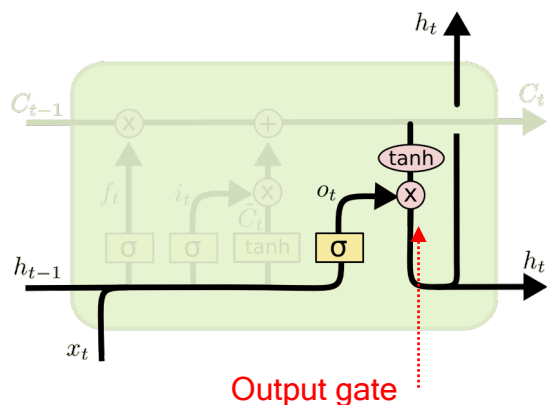
Step 3: Update the cell state



Step 3: Update the cell state $\rightarrow f_t * C_{t-1} + i_t * \tilde{C}_t$

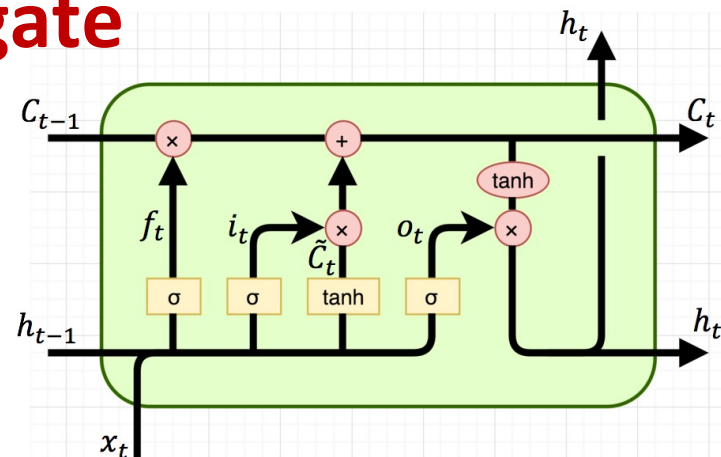
- The new cell state C_t is comprised of information from the past $f_t * C_{t-1}$ and valuable new information $i_t * \tilde{C}_t$
- * denotes elementwise multiplication
- Example: the model drops the old gender information ($f_t * C_{t-1}$) and adds new gender information ($i_t * \tilde{C}_t$) to form the new cell state (C_t)

Step 4: Output gate



$$o_t = \sigma(W_o [h_{t-1}, x_t] + b_o)$$

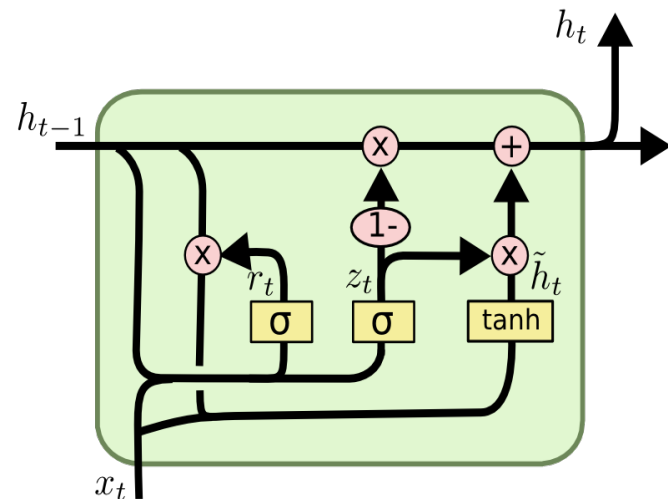
$$h_t = o_t * \tanh(C_t)$$



Step 4: Decide the filtered output from the new cell state $\rightarrow o_t * \tanh(C_t)$

- tanh function filters the new cell state to characterize stored information
 - Significant information in $C_t \rightarrow \pm 1$
 - Minor details $\rightarrow 0$
- Output gate o_t ranges between $[0, 1]$
- h_t serves as a control signal for the next time step
- Example: Since the model just saw a new subject (x_t), it might want to output ($o_t \rightarrow 1$) information relevant to a verb ($\tanh(C_t)$), e.g., singular/plural, in case a verb comes next

Gated Recurrent Unit (GRU)

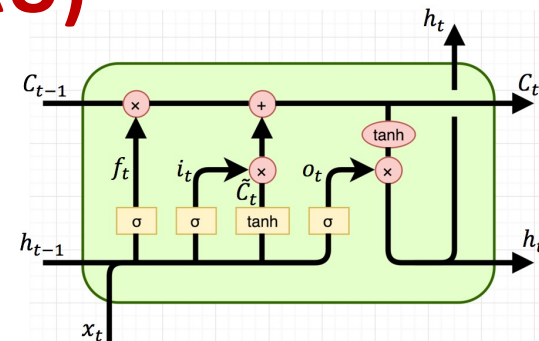


$$z_t = \sigma(W_z \cdot [h_{t-1}, x_t])$$

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

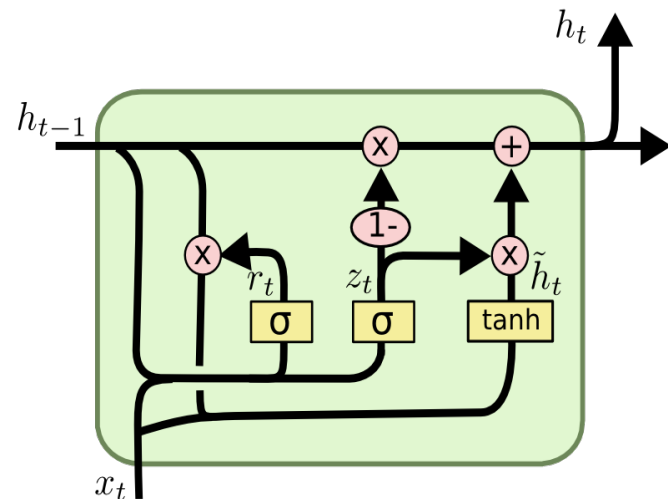
$$\tilde{h}_t = \tanh(W \cdot [r_t * h_{t-1}, x_t])$$

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



- GRU is a variation of LSTM that also adopts the gated design.
- Differences:
 - GRU uses an **update gate** z to substitute the input and forget gates i_t and f_t
 - Combined the cell state C_t and hidden state h_t in LSTM as a single cell state h_t
- GRU obtains similar performance compared to LSTM with fewer parameters and faster convergence. (Cho et al. 2014)

Gated Recurrent Unit (GRU)

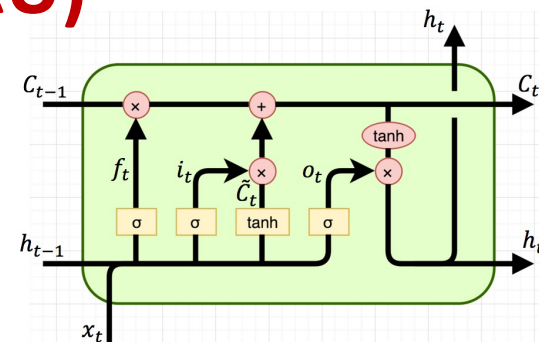


$$z_t = \sigma(W_z \cdot [h_{t-1}, x_t])$$

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

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

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

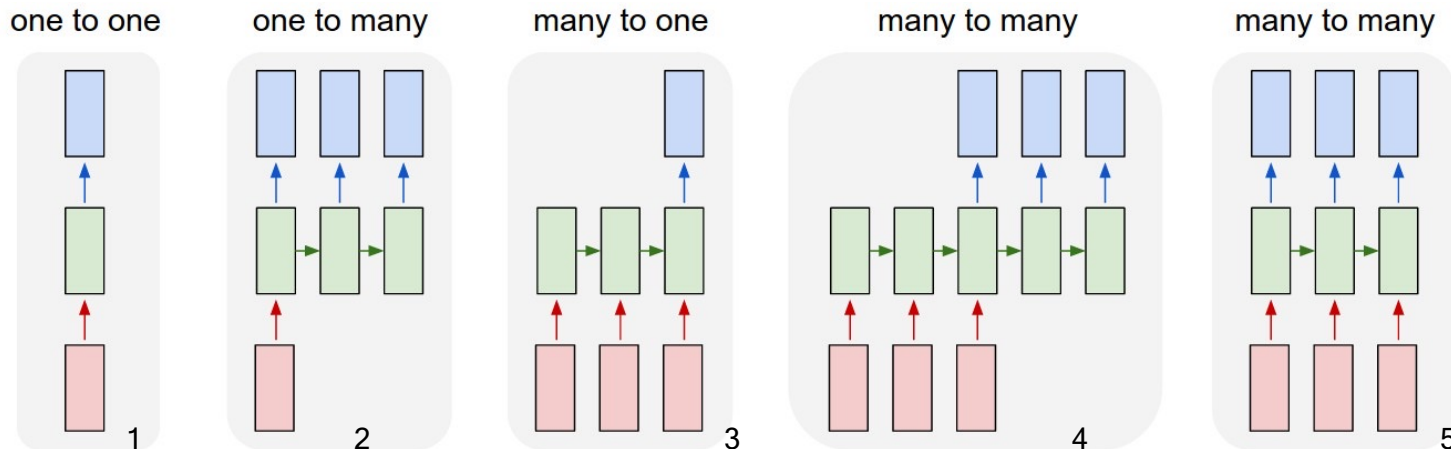


- **Update gate:** controls the composition of the new state
- **Reset gate:** determines how much old information is needed in the alternative state \tilde{h}_t
- **Alternative state:** contains new information
- **New state:** replace selected old information with new information in the new state

Sequence Learning Architectures

- Learning on RNN is more robust when the vanishing/exploding gradient problem is resolved.
 - RNNs can now be applied to different Sequence Learning tasks
- Recurrent NN architecture is flexible to operate over various sequences of vectors.
 - Sequence in the input, the output, or in the most general case both
 - Architecture with one or more RNN layers

Sequence Learning with One RNN Layer

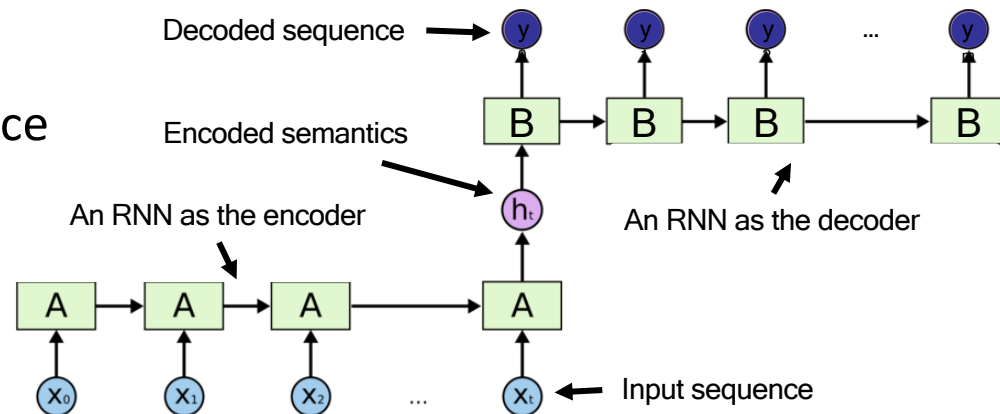


- Each rectangle is a vector and arrows represent functions (e.g., matrix multiply)
- Input vectors are in red, output vectors are in blue and green vectors hold the RNN's state

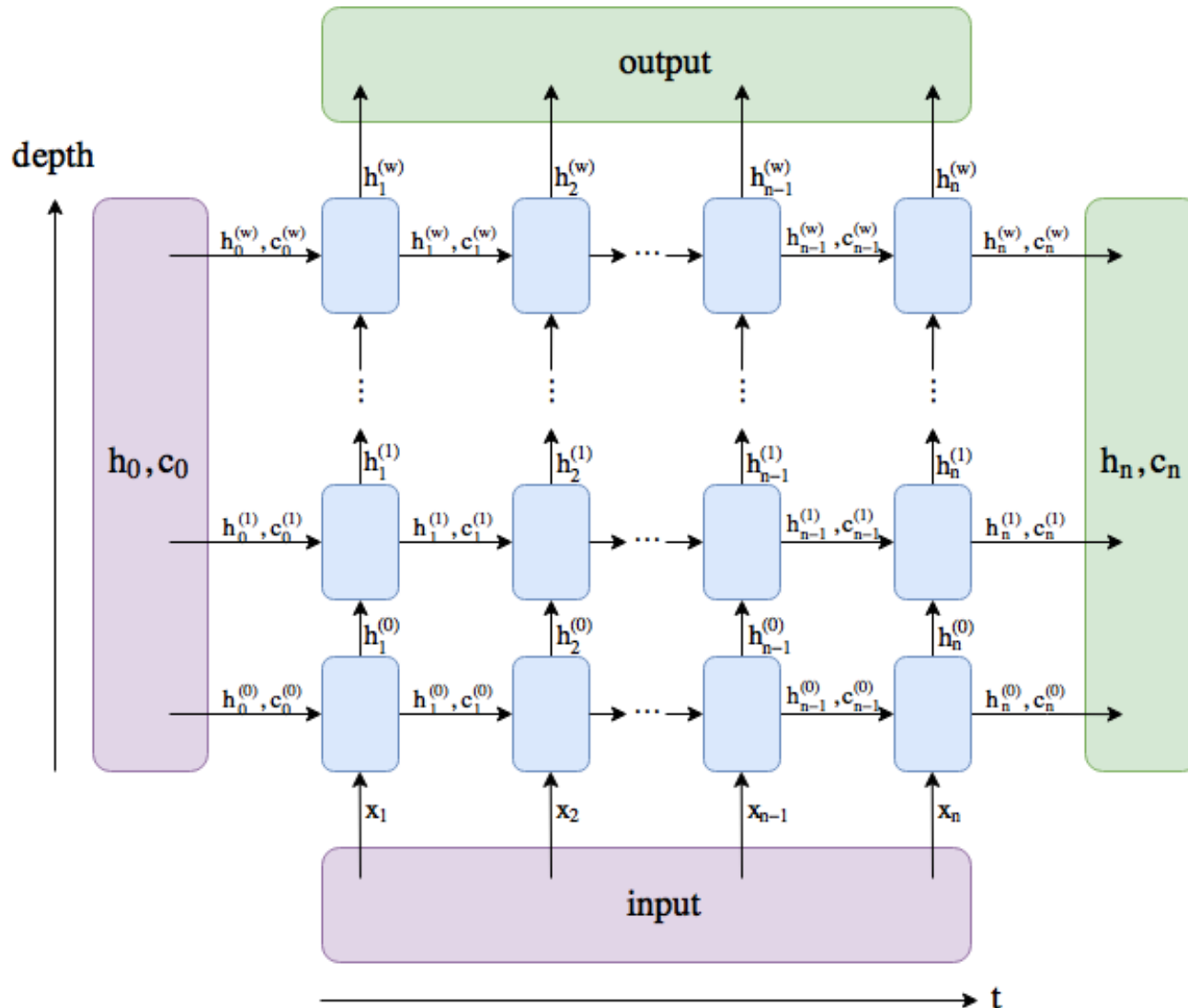
1. Standard NN mode without recurrent structure (e.g., image classification, **one label for one image**)
2. Sequence output (e.g., image captioning, takes an image and outputs a **sentence of words**)
3. Sequence input (e.g., sentiment analysis, a sentence is classified as **expressing positive or negative sentiment**).
4. Sequence input and sequence output (e.g., machine translation, a sentence in **English is translated into a sentence in French**)
5. Synced sequence input and output (e.g., video classification, label each **frame of the video**)

Sequence-to-Sequence (Seq2Seq) model

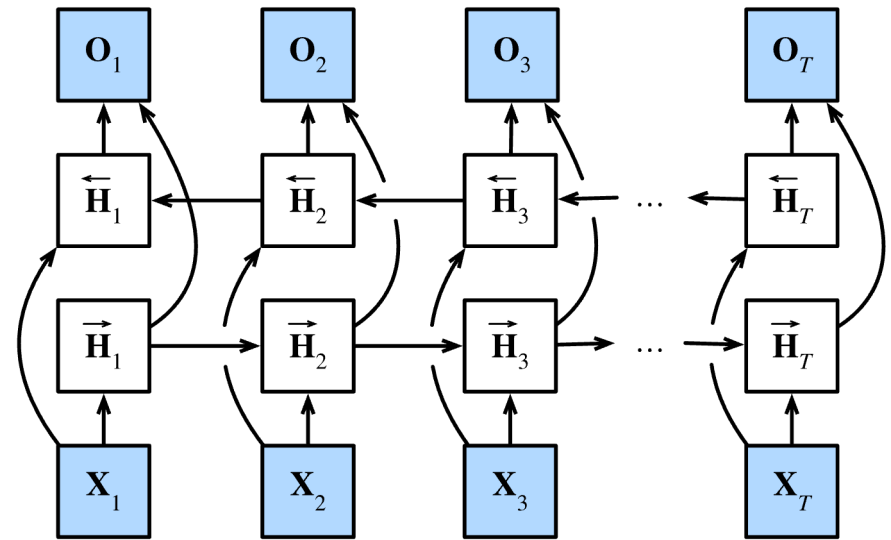
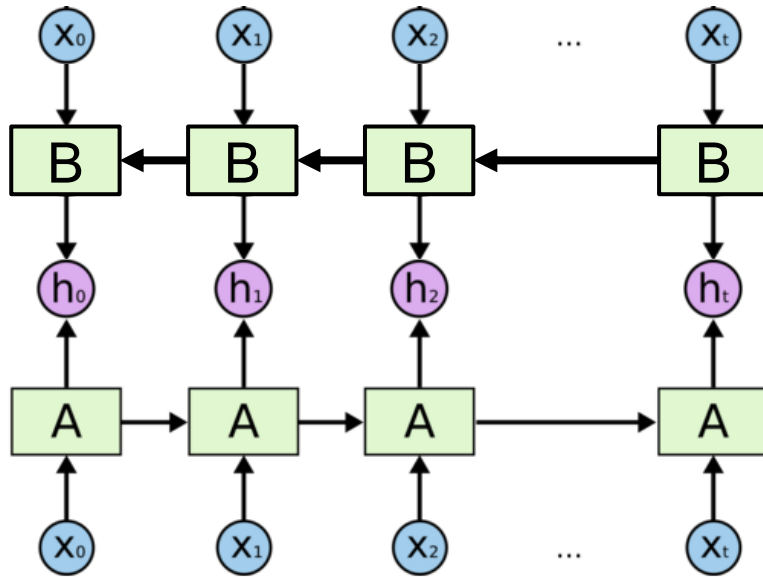
- Developed by Google in 2018 for use in machine translation.
- Seq2seq turns one sequence into another sequence. It does so by use of a [recurrent neural network](#) (RNN) or more often [LSTM](#) or [GRU](#) to avoid the problem of [vanishing gradient](#).
- The primary components are one Encoder and one Decoder network. The encoder turns each item into a corresponding hidden vector containing the item and its context. The decoder reverses the process, turning the vector into an output item, using the previous output as the input context.
- **Encoder RNN:** extract and compress the semantics from the input sequence
- **Decoder RNN:** generate a sequence based on the input semantics
- Apply to tasks such as machine translation
 - Similar underlying semantics
 - E.g., “I love you.” to “Je t’aime.”



num_layers



Bidirectional RNN

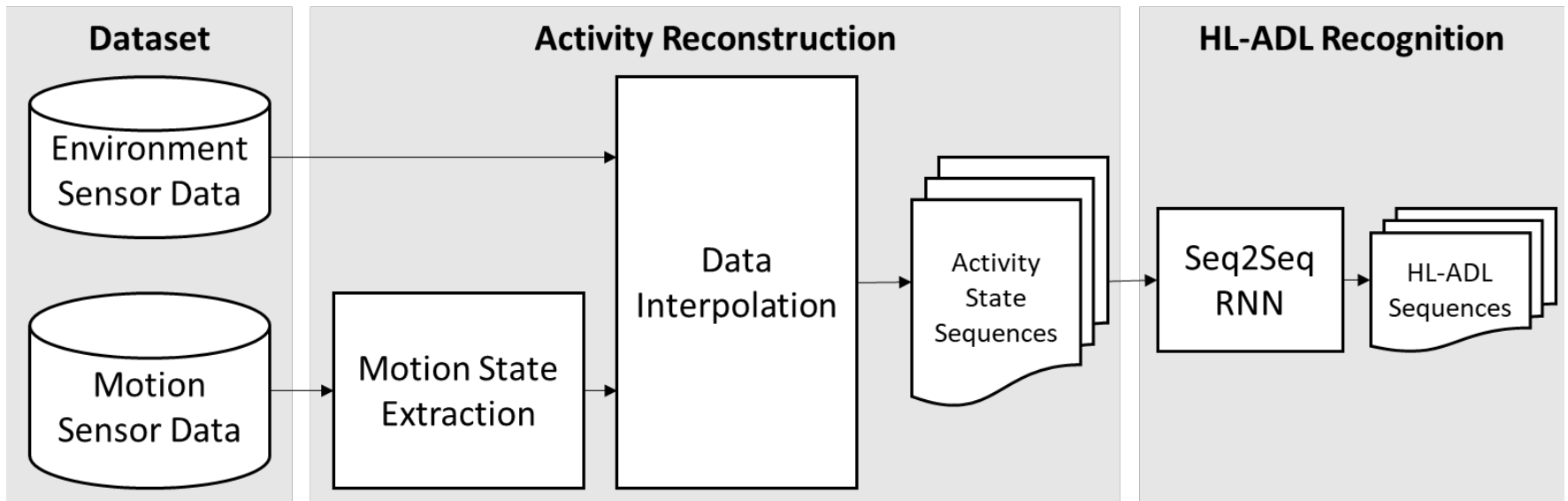


- Connects two recurrent units (synced many-to-many model) of opposite directions to the same output.
- Captures forward and backward information from the input sequence
- Apply to data whose current state (e.g., h_0) can be better determined when given future information (e.g., x_1, x_2, \dots, x_t)
 - E.g., in the sentence “the bank is robbed,” the semantics of “bank” can be determined given the verb “robbed.”

RNN Application: High-level ADL Recognition

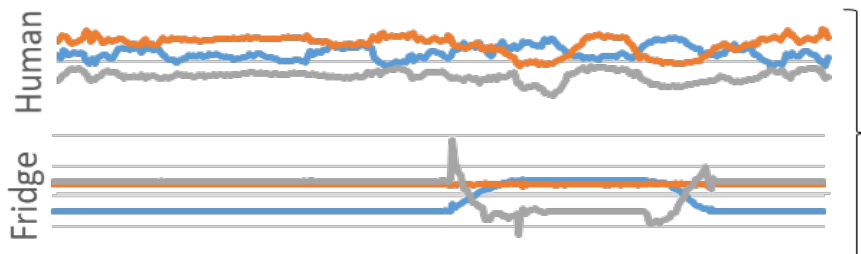
- Activities of Daily Living (ADLs) are introduced to evaluate the self-care ability of senior citizens (Williams 2014)
- Sensor-based home monitoring systems
 - **Environment:** Placed in the environment to capture changes (e.g., on/off, motion)
 - **Wearable:** Attached to the body to measure movements and physiological signals
- Sensors sample at 10 Hz → 0.8 million records per day
- Machine learning (esp. deep learning) needed to recognize ADLs from the large amount of data
 - Mid-level (ML) ADLs: gestures, gait, etc.
 - **High-level (HL) ADLs: preparing food, medical intake, senior care, etc.**
- Research Objective: develop a universal ADL recognition framework to extract HL-ADLs from raw sensor data

Seq2Seq-ADL Research Design



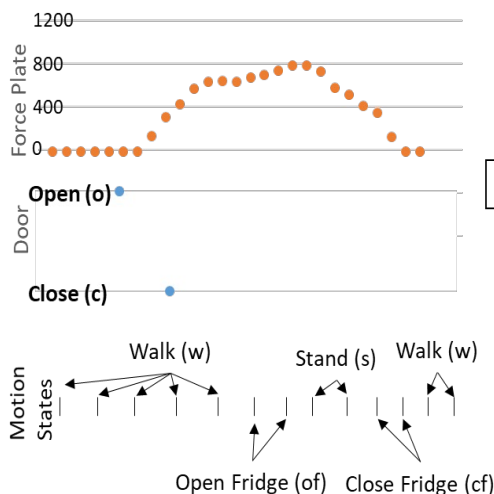
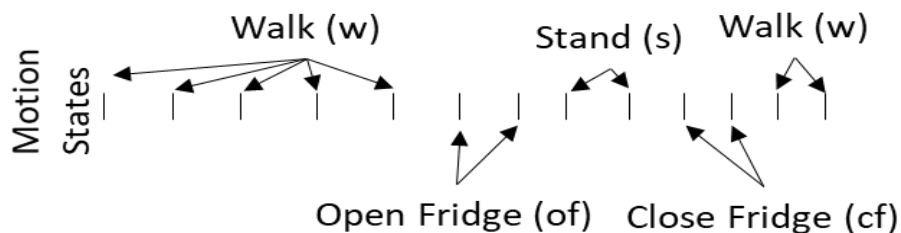
- Intuition: Recognizing the HL-ADLs from the sensor data is similar to captioning/translation. We can generate ADL label sequence with a Seq2Seq model for the input data. The underlying semantics are similar.
- A Seq2Seq model is designed to extract HL-ADLs from the activity state sequence

Activity Reconstruction



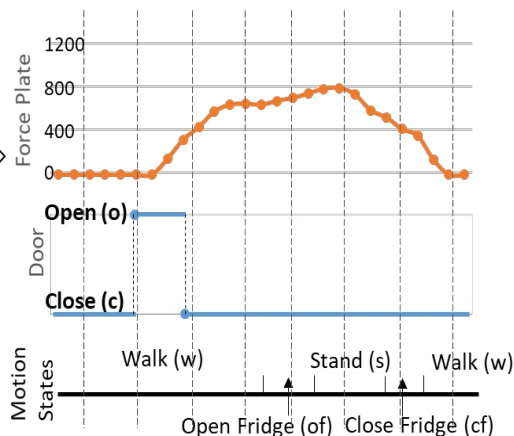
Step 1:

Motion State
Extraction



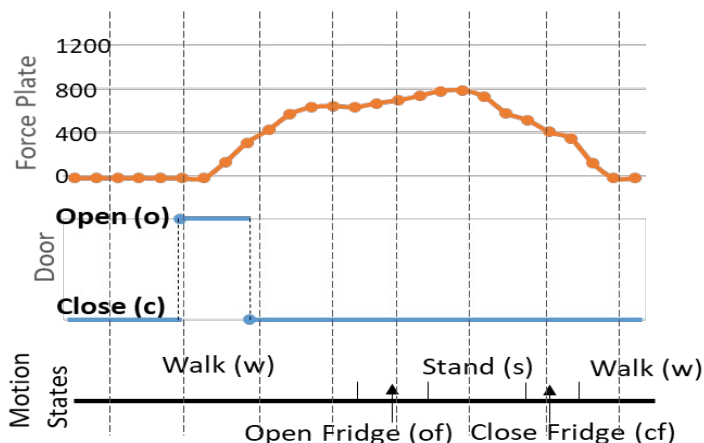
Step 2:

Data Interpolation



- Objective: create temporally aligned activity representation from different data sources
- Four sensors for demonstration:
 - a force plate → pressure on the floor
 - a door on/off sensor → open (o) and close (c) states
 - a human motion sensor
 - object motion sensor attached on the fridge
- **Step 1.** Extract discrete motion states from motion sensor data with a state-of-the-art gesture recognition model
- **Step 2.** Interpolate discrete data from each sensor

Activity Reconstruction



Step 3:

Sample
Activity States

$$S = \begin{bmatrix} E_1 \\ E_2 \\ M \end{bmatrix} = \begin{bmatrix} 0 \\ c \\ w \end{bmatrix} \begin{bmatrix} 0 \\ o \\ w \end{bmatrix} \begin{bmatrix} 390 \\ c \\ w \end{bmatrix} \begin{bmatrix} 633 \\ c \\ w \end{bmatrix} \begin{bmatrix} 710 \\ c \\ of \end{bmatrix} \begin{bmatrix} 794 \\ c \\ s \end{bmatrix} \begin{bmatrix} 400 \\ c \\ cf \end{bmatrix} \begin{bmatrix} 0 \\ c \\ w \end{bmatrix}$$

Step 4:

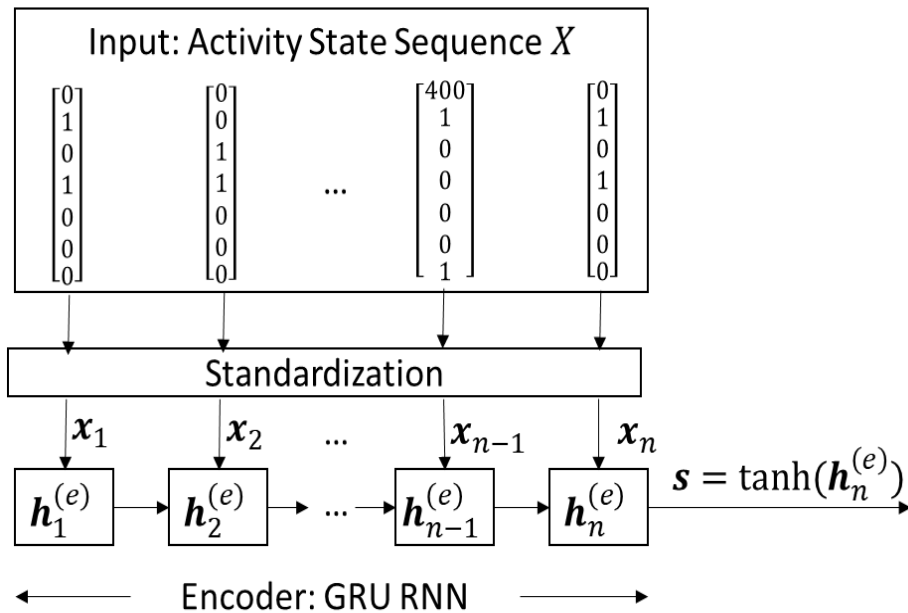
One-hot
encoding

$$S = \begin{bmatrix} E_1 \\ E_2 = c \\ E_2 = o \\ M = w \\ M = of \\ M = s \\ M = cf \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 1 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 1 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

Step 5: $X = [S_1 \ S_2 \ S_3 \ S_4 \ S_5 \ S_6 \ S_7 \ S_8]$

- **Steps 3.** Sample each data stream at same timestamps to construct the Activity State representations
 - Temporally aligned observations
- **Steps 4.** Encode the Activity States S_i
 - Encode categorical (discrete) values using one-hot encoding
- **Step 5.** Organize the states vector in a matrix X
 - Data matrix X aggregates temporally aligned sensor observation sequences to represent the activity

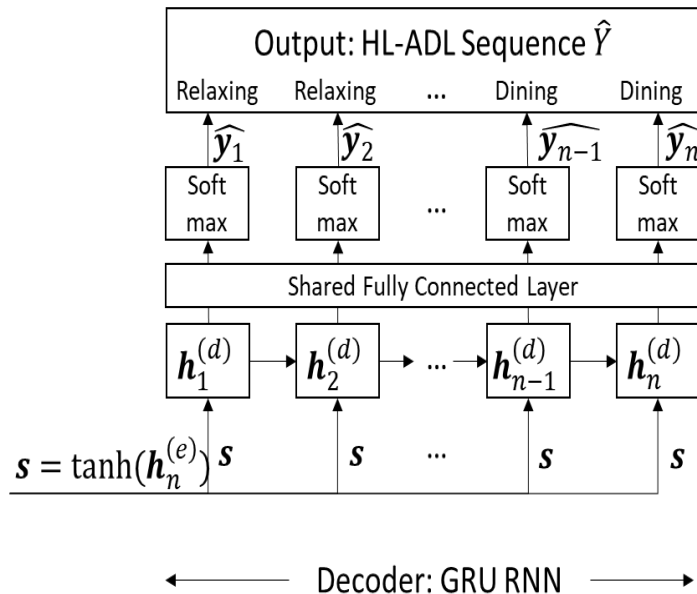
Seq2Seq - Encoder



- The encoder network takes the Activity State Sequence X as the input to generate the activity semantics vector s
- The encoder network adopts GRU recurrent cells to learn temporal patterns
 - Each hidden state $h_i^{(e)}$ depends on the input x_i and the previous state $h_{i-1}^{(e)}$

- s is expected to be a condensed representation for human/object motions, object usages, and their temporal patterns during the period.

Seq2Seq - Decoder



- The decoder network takes the encoded activity semantics vector s to generate HL-ADL label for each input vector
- The decoder network also adopts GRU recurrent cells to interpret the temporal patterns
- Multi-class classification
 - Softmax \rightarrow probability distribution over output classes
 - Categorical cross-entropy loss

Experimental results

Recall of S2S_GRU against benchmarks on a wearable/environment motion sensor dataset

Recall (%)	Relaxing	Coffee time	Early morning	Cleanup	Sandwich time
S2S_GRU	61.7	66.5	72.1	82.6	79.1
HMM	48.8	19.5***	6.4***	17.4***	40.9***
S2S_LSTM	60.0	54.9**	75.9	83.2	59.9***

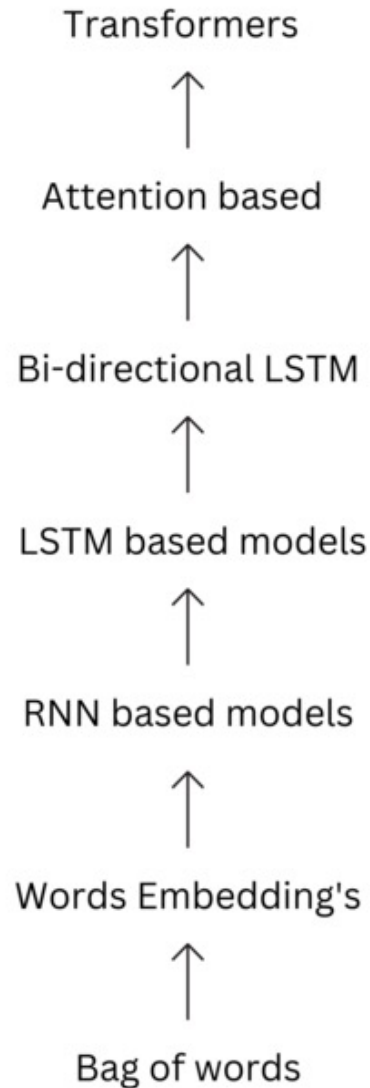
Recall of S2S_GRU against benchmarks on an environment sensor dataset

Recall (%)	Filling pillbox	Watch DVD	Water plants	Answer the phone	Prepare gift card	Prepare soup	Cleaning	Choose outfit
S2S_GRU	88.6	79.0	58.5	80.8	92.1	71.1	85.0	58.0
HMM	33.2***	31.6***	16.8***	22.0***	29.9***	26.4***	20.4***	22.5***
S2S_LSTM	86.1*	77.8	64.3	74.9*	90.8*	69.6	82.5	49.5**

*: p-value<0.05, **: p-value<0.01, ***: p-value<0.001

- S2S_GRU model is evaluated on two different datasets
- S2S_GRU is more accurate and flexible in adjusting to different real-life HL-ADL recognition applications

Development of natural language processing tools



Summary

- LSTM and GRU are RNNs that retain past information and update with a gated design.
 - The “additive” structure avoids vanishing gradient problem
- RNNs allow flexible architecture designs to adapt to different sequence learning requirements.
- RNNs have broad real-life applications.
 - Text processing, machine translation, signal extraction/recognition, image captioning
 - Mobile health analytics, activity of daily living, senior care

Important References

- Bengio, Y., Simard, P., & Frasconi, P. (1994). Learning long-term dependencies with gradient descent is difficult. *IEEE transactions on neural networks*, 5(2), 157-166.
- Goodfellow, I., Bengio, Y., Courville, A., & Bengio, Y. (2016). *Deep learning* (Vol. 1). Cambridge: MIT press.
- Graves, A. (2012). *Supervised sequence labelling with recurrent neural networks*. <https://www.cs.toronto.edu/~graves/preprint.pdf>
- Jozefowicz, R., Zaremba, W., & Sutskever, I. (2015, June). An empirical exploration of recurrent network architectures. In *International Conference on Machine Learning* (pp. 2342-2350).
- Lipton, Z. C., Berkowitz, J., & Elkan, C. (2015). A critical review of recurrent neural networks for sequence learning. *arXiv preprint arXiv:1506.00019*.
- Salehinejad, H., Baarbe, J., Sankar, S., Barfett, J., Colak, E., & Valaee, S. (2017). Recent Advances in Recurrent Neural Networks. *arXiv preprint arXiv:1801.01078*.