



Birla Institute of Technology and Science Pilani, Hyderabad Campus  
2<sup>nd</sup> Semester 2023-24

16.04.2024

# **BITS F464: Machine Learning**

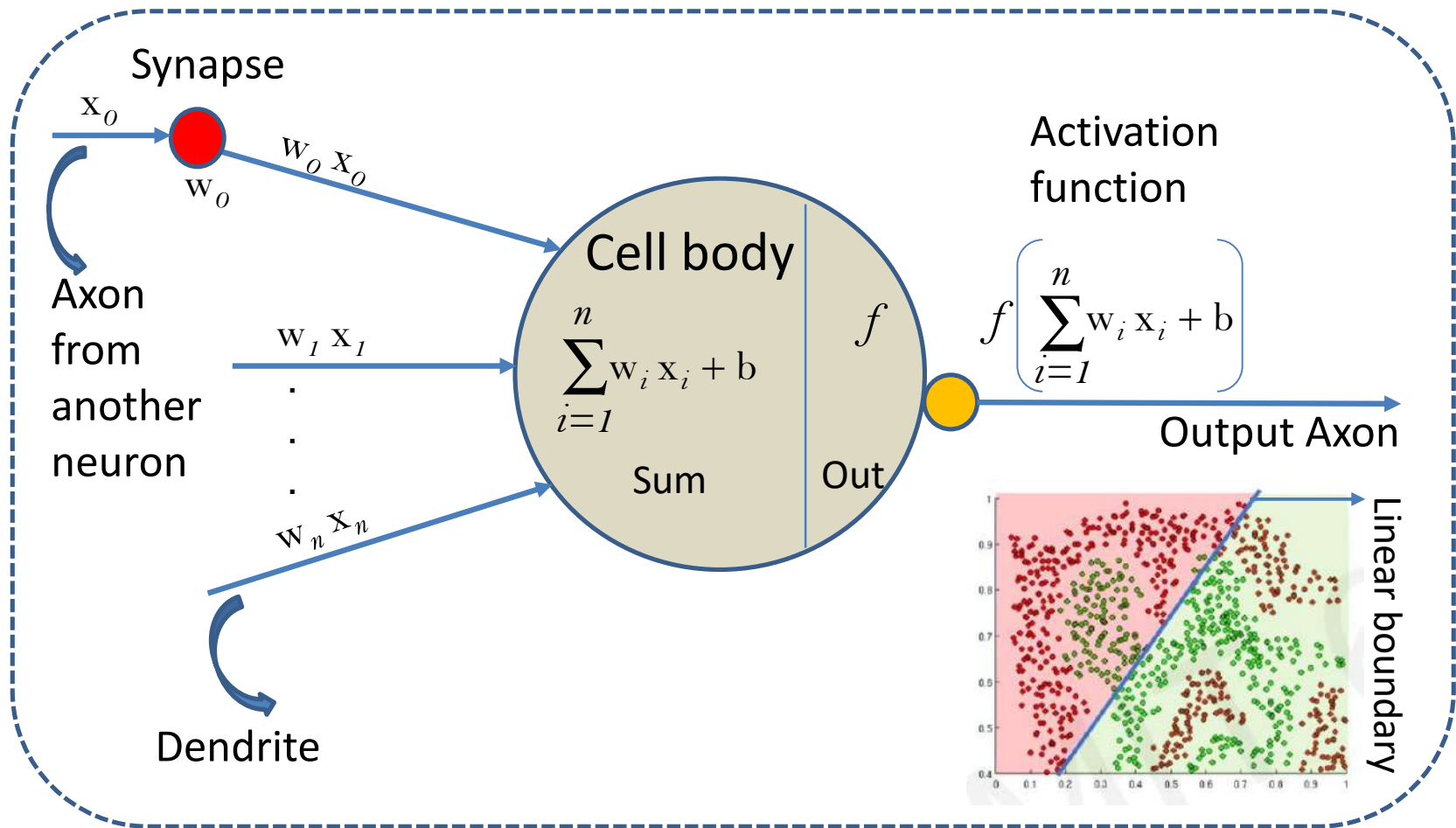
---

**NEURAL NETWORKS: CONVOLUTIONAL/ RECURRENT/ GENERATIVE MODELS**

Chittaranjan Hota, Sr. Professor  
Dept. of Computer Sc. and Information Systems  
[hota@hyderabad.bits-pilani.ac.in](mailto:hota@hyderabad.bits-pilani.ac.in)

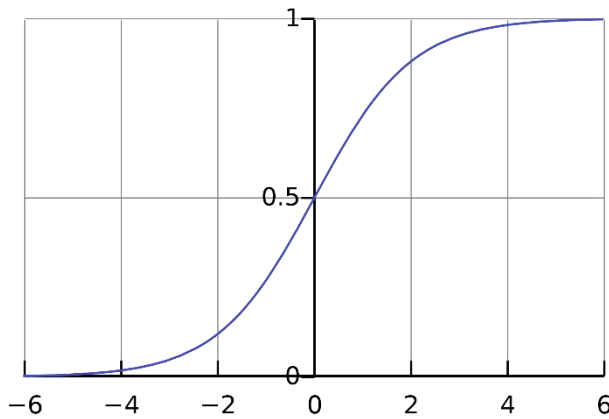
---

# Recap: A Perceptron



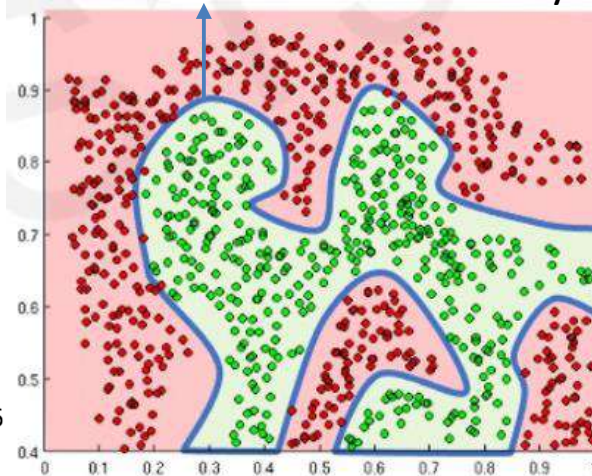
# Vanishing Gradient in MLPs/BPNs

$$\sigma(z) = 1/(1+e^{-z})$$



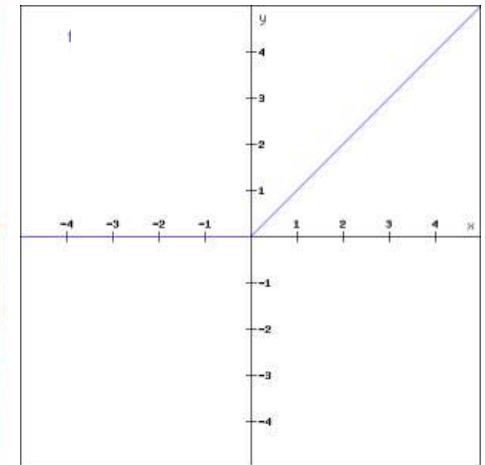
Sigmoid function/ Logistic Curve

Non-linear boundary



Approx. Arbitrary Complex fun

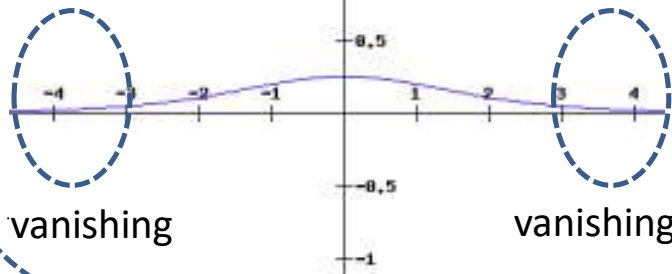
$$\text{ReLU}(z) = \max(0, z)$$



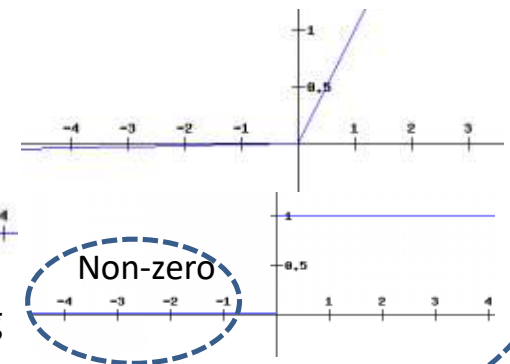
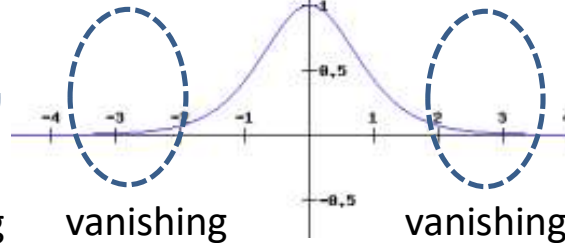
$$\text{Leaky ReLU}(z) = 0.01z, z < 0 \\ = z, z \geq 0$$

$$\sigma(z) \cdot (1 - \sigma(z))$$

derivative



$$1/\cosh^2(z)$$

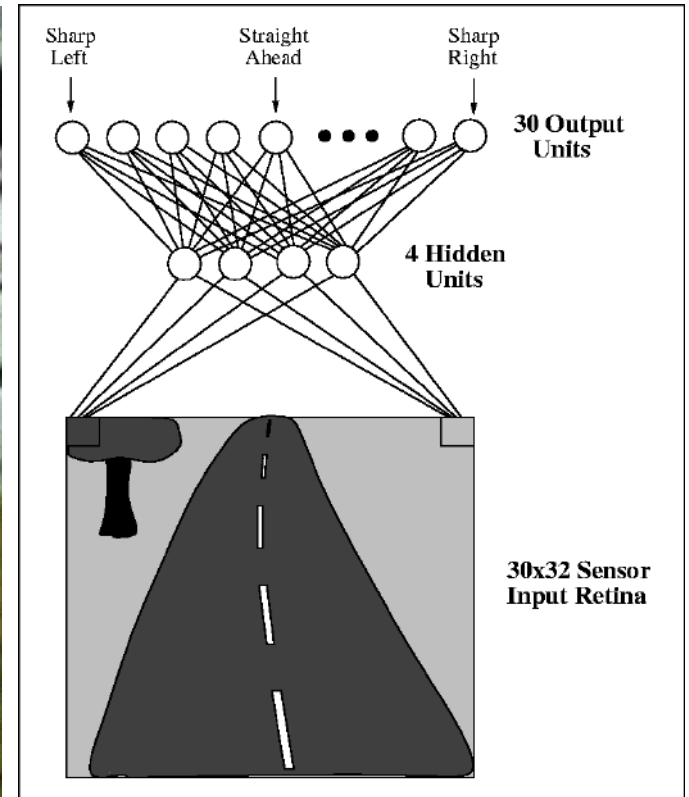


Non-zero



# ALVINN: An Autonomous Land Vehicle In a Neural Network

---



Source: <https://www.ri.cmu.edu/>

(1989: 3-layer Network)

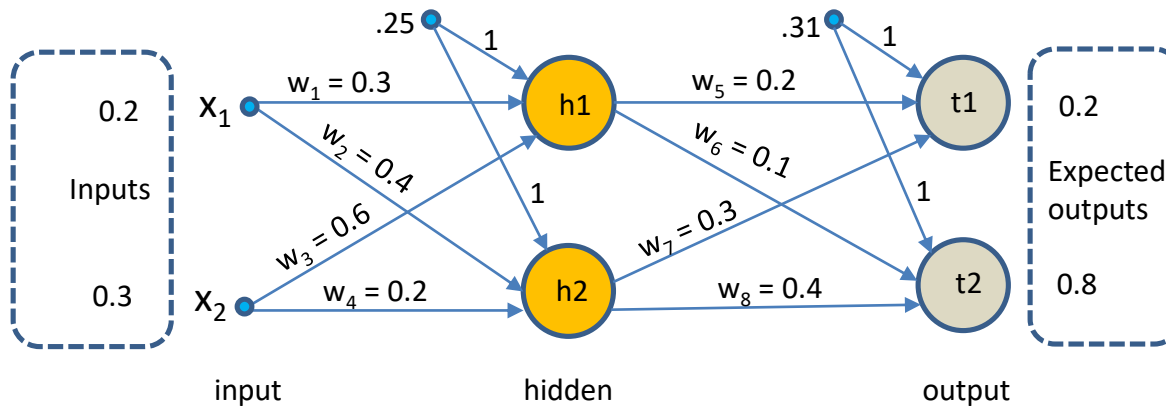
---

An application of a Backpropagation Neural Network in smart driving

# Example Backpropagation Neural Networks

## Quiz Question

Learning rate:  $\eta = 0.4$



Total error:

$$E_{\text{total}} = \frac{1}{2} \sum (\text{target} - \text{output})^2$$

$$E_1 = \frac{1}{2} (0.2 - .6486)^2 = .1006$$

$$E_2 = \frac{1}{2} (0.8 - .6480)^2 = .0116$$

## Forward Pass:

For h1:

$$\text{Sum} = .25 \times 1 + .3 \times .2 + .3 \times .6 = 0.49$$

$$\text{Output} = \frac{1}{1 + e^{-.49}} = \mathbf{0.6201}$$

For h2:

$$\text{Sum} = .25 \times 1 + .2 \times .4 + .3 \times .2 = 0.39$$

$$\text{Output} = \frac{1}{1 + e^{-.39}} = \mathbf{0.5963}$$

For t1:

$$\text{Sum} = .31 \times 1 + .6201 \times .2 + .5963 \times .3 = .6129$$

$$\text{Output} = \frac{1}{1 + e^{-.6129}} = \mathbf{0.6486}$$

For t2:

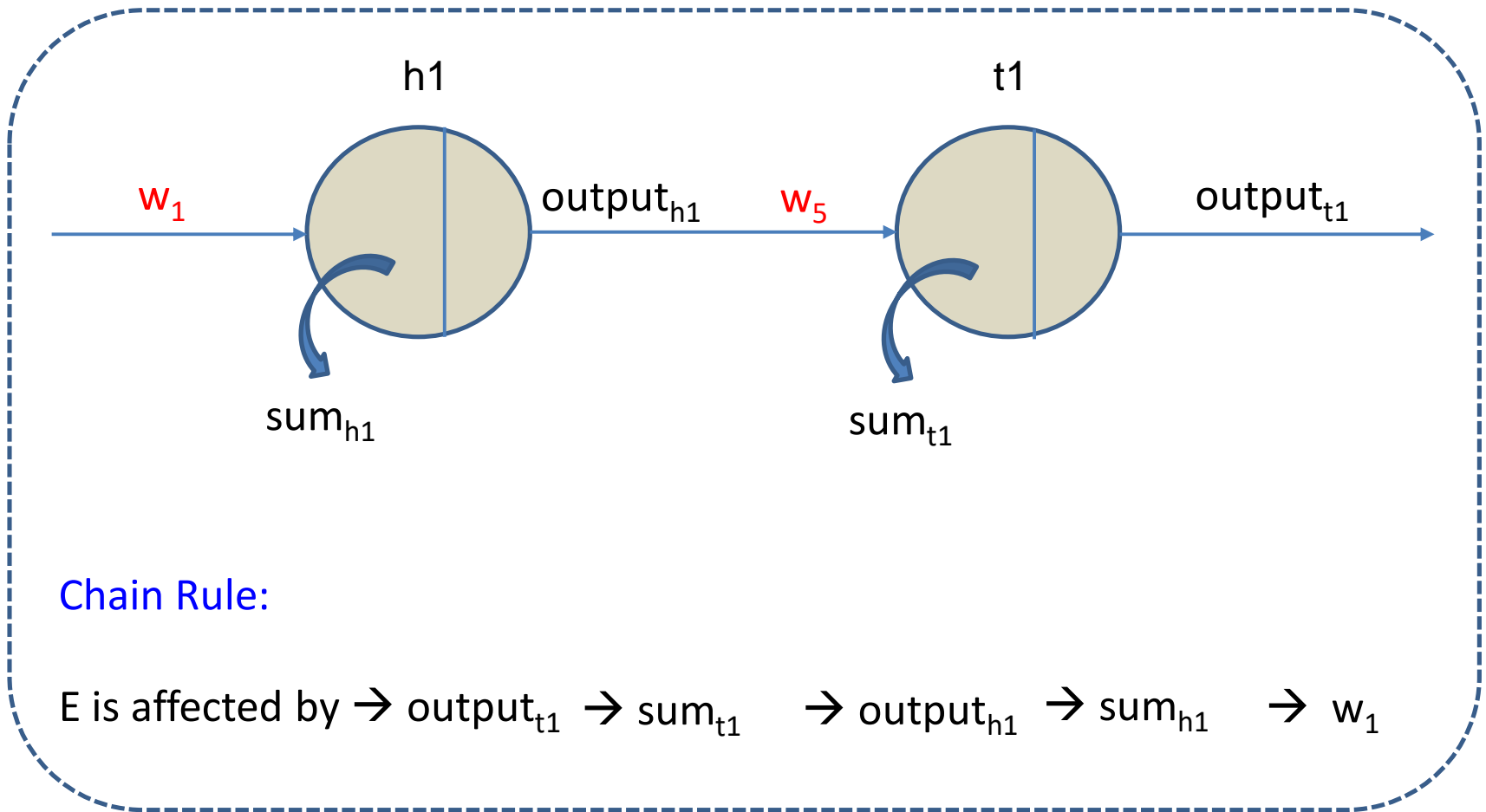
$$\text{Sum} = .31 \times 1 + .6201 \times .1 + .5963 \times .4 = .6105$$

$$\text{Output} = \frac{1}{1 + e^{-.6105}} = \mathbf{0.6480}$$

$$E_{\text{total}} = .1006 + .0116 = 0.1122$$

# Example Chain Rule

---



# Backward Pass

$$\frac{\partial E_{\text{total}}}{\partial w_5} = \frac{\partial E_{\text{total}}}{\partial \text{output}_{t1}} \times \frac{\partial \text{output}_{t1}}{\partial \text{sum}_{t1}} \times \frac{\partial \text{sum}_{t1}}{\partial w_5}$$

$$E_{\text{total}} = \frac{1}{2} \sum (\text{actual} - \text{observed})^2 \quad \dots \text{Eq. (1)}$$

$$= \frac{1}{2} \left\{ (\text{actual}_{t1} - \text{output}_{t1})^2 + (\text{actual}_{t2} - \text{output}_{t2})^2 \right\}$$

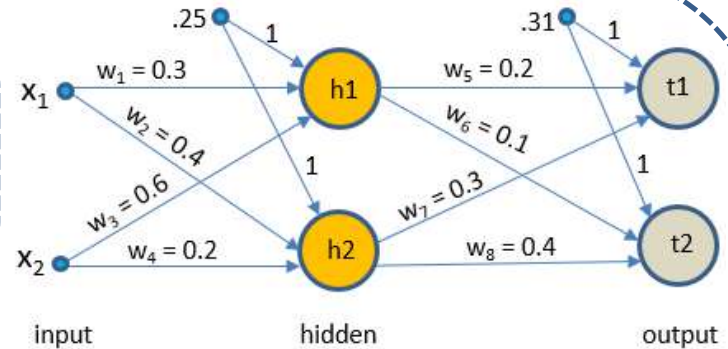
$$\frac{\partial E_{\text{total}}}{\partial \text{output}_{t1}} = 2 \times \frac{1}{2} (\text{actual}_{t1} - \text{output}_{t1}) \times -1 + 0$$

$$= (\text{output}_{t1} - \text{actual}_{t1})$$

$$= 0.6486 - 0.2 = \mathbf{0.4486} \quad \dots \text{Eq. (2)}$$

$$\frac{\partial \text{sum}_{t1}}{\partial w_5} = \frac{\partial (\text{output}_{h1} \times w_5 + \text{output}_{h2} \times w_7)}{\partial w_5}$$

$$= \text{output}_{h1} = \mathbf{0.6201} \quad \dots \text{Eq. (4)}$$



$$\frac{\partial \text{output}_{t1}}{\partial \text{sum}_{t1}} = ?$$

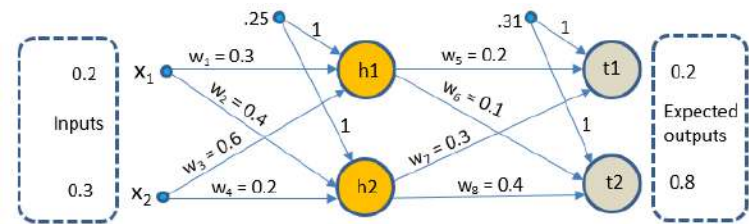
$$\sigma(x) = 1 / (1 + e^{-x})$$

$$d/dx(\sigma(x)) = \sigma(x) (1 - \sigma(x))$$

$$\frac{\partial \text{output}_{t1}}{\partial \text{sum}_{t1}} = \text{output}_{t1} (1 - \text{output}_{t1})$$

$$= 0.6486 (1 - 0.6486) = \mathbf{0.2279} \quad \dots \text{Eq. (3)}$$

# Continued...



$$\text{Eq. (1): } \frac{\partial E_{\text{total}}}{\partial w_5} = \text{Eq.(1)} \times \text{Eq.(2)} \times \text{Eq.(3)} = .4486 \times .2279 \times .6201 = .0634 \rightarrow$$

$$w_5 = w_5 - \eta \frac{\partial E_{\text{total}}}{\partial w_5} = .2 - .4 \times .0634 = .1746$$

$$\frac{\partial E_{\text{total}}}{\partial w_6} = \frac{\partial E_{\text{total}}}{\partial \text{output}_{t_2}} \times \frac{\partial \text{output}_{t_2}}{\partial \text{sum}_{t_2}} \times \frac{\partial \text{sum}_{t_2}}{\partial w_6} = (.6480 - .8) \times (.6480 \times (1 - .6480)) \times .6201$$

$$= -.152 \times .2281 \times .6201 = -.0215 \rightarrow w_6 = w_6 - \eta \frac{\partial E_{\text{total}}}{\partial w_6} = .1 - (.4 \times -.0215) = .1086$$

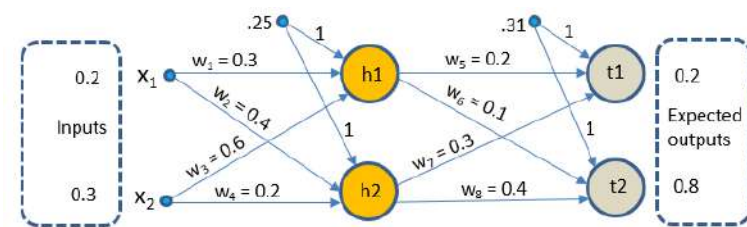
$$\frac{\partial E_{\text{total}}}{\partial w_7} = \frac{\partial E_{\text{total}}}{\partial \text{output}_{t_1}} \times \frac{\partial \text{output}_{t_1}}{\partial \text{sum}_{t_1}} \times \frac{\partial \text{sum}_{t_1}}{\partial w_7} = \text{Eq.(2)} \times \text{Eq.(3)} \times \text{output}_{h_2}$$

$$= 0.4486 \times 0.2279 \times 0.5963 = 0.0609$$

$$\rightarrow w_7 = w_7 - \eta \frac{\partial E_{\text{total}}}{\partial w_7} = .3 - .4 \times .0609 = 0.3 - 0.02436 = 0.2756$$



# Continued...



$$\frac{\partial E_{\text{total}}}{\partial w_8} = \frac{\partial E_{\text{total}}}{\partial \text{output}_{t_2}} \times \frac{\partial \text{output}_{t_2}}{\partial \text{sum}_{t_2}} \times \frac{\partial \text{sum}_{t_2}}{\partial w_8} = (-.152) \times .2281 \times \text{output}_{h_2} = -0.0207$$

$$\rightarrow w_8 = w_8 - \eta \frac{\partial E_{\text{total}}}{\partial w_8} = 0.4 + 0.4 \times 0.0207 = 0.4 + 0.0083 = .4083$$

Similarly,  
 $w_2 = .4007$   
 ... for you

Now Compute Weights in the Hidden Layer ( $w_1, w_2, w_3,$  and  $w_4$ ): Chain becomes longer.

For  $w_1$ :

$$\frac{\partial E_{\text{total}}}{\partial w_1} = \frac{\partial E_1}{\partial w_1} + \frac{\partial E_2}{\partial w_1}$$

$$w_1 = w_1 + \eta \left( \frac{\partial E_{\text{total}}}{\partial w_1} \right) = 0.3 - 0.4 \times 0.0008 = .2997$$

Where,

$$\frac{\partial E_1}{\partial w_1} = \frac{\partial E_1}{\partial \text{output}_{t_1}} \times \frac{\partial \text{output}_{t_1}}{\partial \text{sum}_{t_1}} \times \frac{\partial \text{sum}_{t_1}}{\partial \text{output}_{h_1}} \times \frac{\partial \text{output}_{h_1}}{\partial \text{sum}_{h_1}} \times \frac{\partial \text{sum}_{h_1}}{\partial w_1}$$

$$\rightarrow \frac{\partial E_1}{\partial w_1} = .4486 \times .2279 \times w_5 \times (\text{output}_{h_1} \times (1 - \text{output}_{h_1})) \times 0.2 = 0.00096$$

Now,

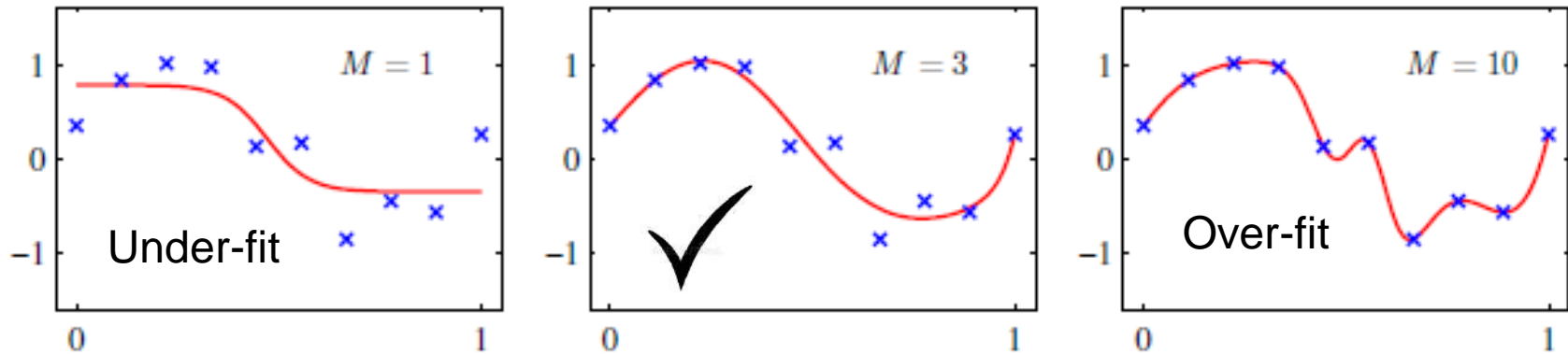
$$\frac{\partial E_2}{\partial w_1} = \frac{\partial E_2}{\partial \text{output}_{t_2}} \times \frac{\partial \text{output}_{t_2}}{\partial \text{sum}_{t_2}} \times \frac{\partial \text{sum}_{t_2}}{\partial \text{output}_{h_1}} \times \frac{\partial \text{output}_{h_1}}{\partial \text{sum}_{h_1}} \times \frac{\partial \text{sum}_{h_1}}{\partial w_1}$$

$$= -.1520 \times .2281 \times w_6 \times .2356 \times .2 = -.00016 \rightarrow \frac{\partial E_{\text{total}}}{\partial w_1} = .00096 - .00016 = .0008$$

# Regularization in Neural Networks

---

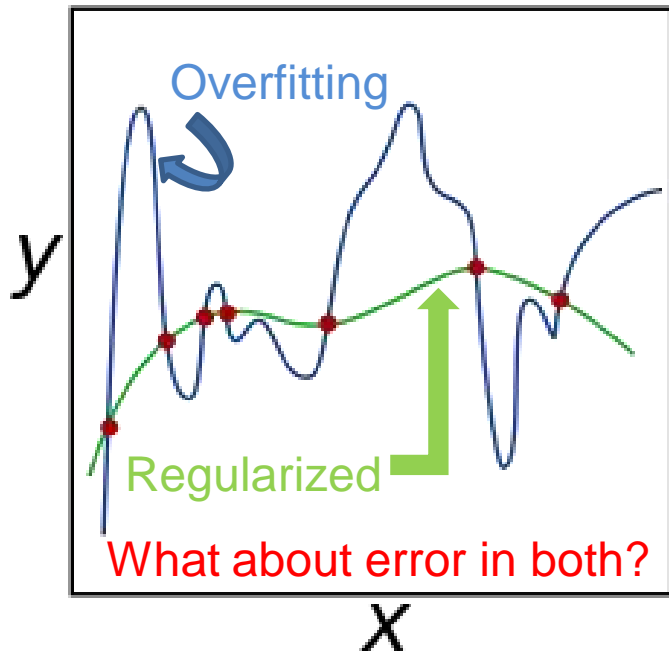
- Which one is a free parameter in a Neural network?
  - Input ~~X~~ Output or Number of units in the hidden layer ( $M$ )
- Why Regularization is needed in Neural Networks?
  - To improve the generalization/ learning outcome. To control impact of noise and fluctuations on the dataset. Alternatively, to avoid over-fitting.



(Fitting a Sinusoidal dataset with different number of hidden units and **Sum-of-Squares** error function optimized by Gradient descent)

# Regularization: **Weight Decay (L1/L2)**

- Control model complexity by the addition of a regularization term to the error function.



Loss function:  $E_D(\mathbf{w}) + \lambda E_W(\mathbf{w})$

Data term:  $E_D(\mathbf{w})$

Regularization term:  $\lambda E_W(\mathbf{w})$

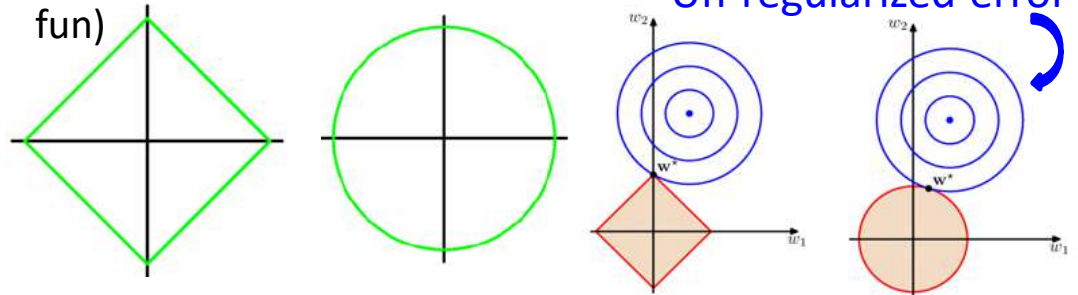
With Sum-of-squares error and a quadratic regularizer:

$$\frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\}^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}$$

Contour plots  
(behaviour of  
fun)

L2 regularization (**Ridge**)

Un-regularized error

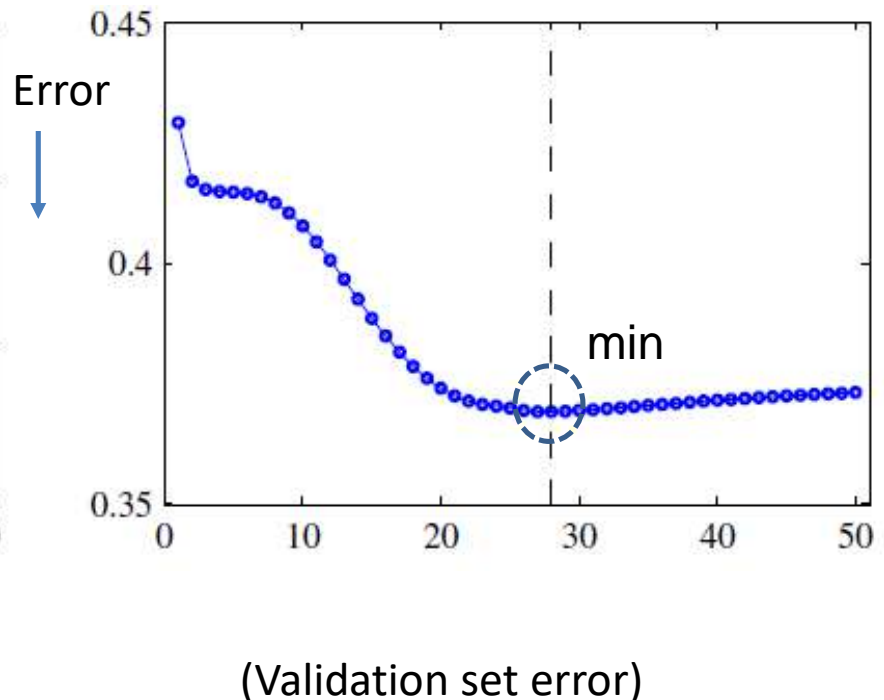
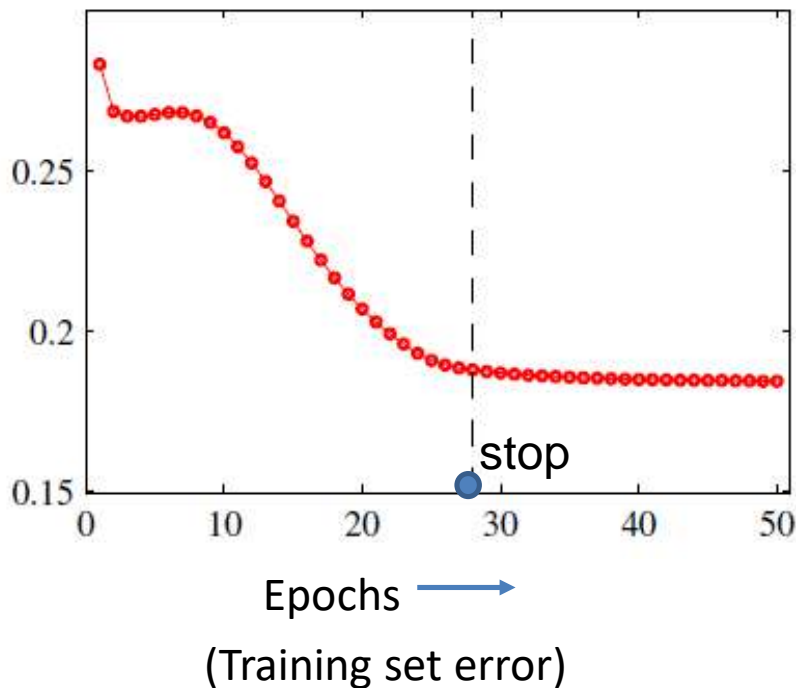


In L1 regularization (**Lasso**), the additional term added to the loss function is the sum of the absolute values of the weights. This encourages sparsity in the weights, effectively **shrinking some of them to zero**.  $J_{L1}(\theta) = J(\theta) + \lambda \sum_{i=1}^n |\theta_i|$

Where,  $\lambda$  as the regularization parameter,  $\vartheta$  as the vector of weights of the Network.

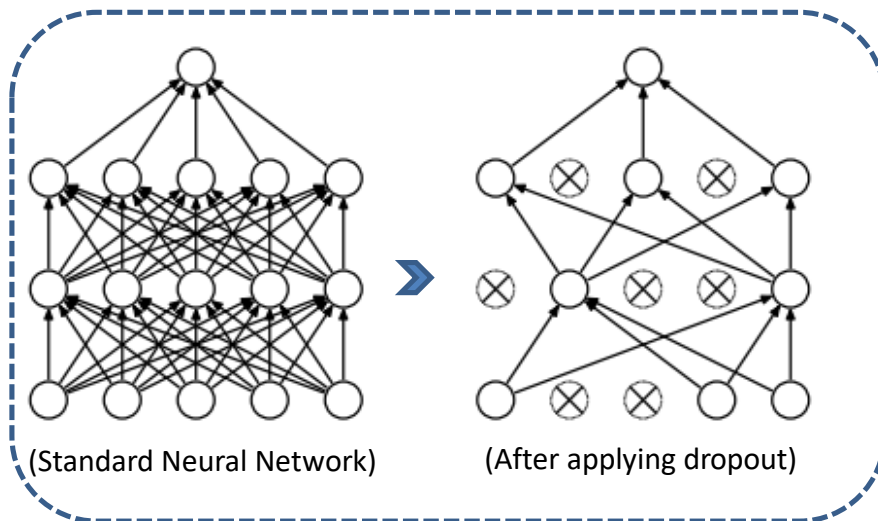
# Regularization: Early Stopping

- Early stopping monitors the performance of the model on a validation set and stops training when the performance starts to degrade, thus preventing the model from overfitting to the training data.
- Example sinusoidal dataset

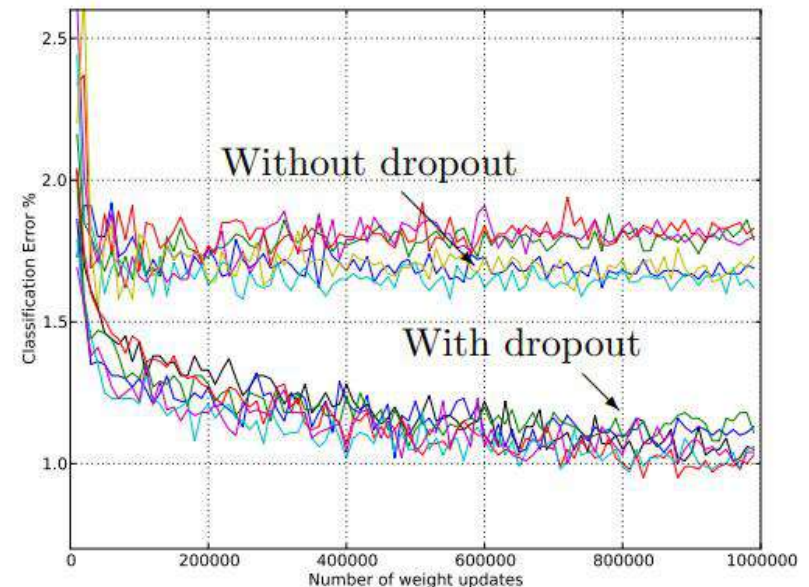


# Stochastic Regularization: Dropout

- Drop out each individual unit with some probability  $p$  (usually  $p = 1/2$ ) by setting its activation to '0'.
- The key idea behind dropout is to prevent overfitting by adding noise to the network during training.



(Img. Source: Nitish Srivastava, et al., Journal of Machine Learning Research, 15, 2014)



During inference (testing or prediction), dropout is typically turned off, and the full network is used. However, the weights are **usually scaled** by '1-p' during inference to account for the fact that more units were active during training.



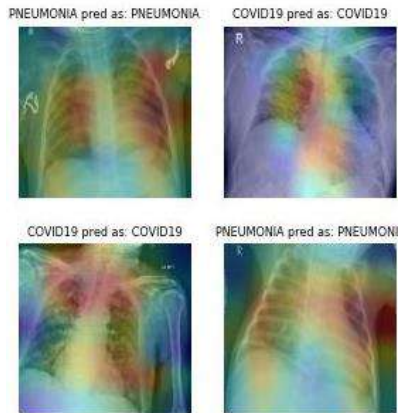
# Regularization: Data Augmentation

- Data augmentation acts as a form of regularization by introducing additional **variations** and **diversity** into the training dataset.
- A technique used to artificially increase the size of a training dataset by applying various transformations to the existing data samples.
- Random rotation, Random scaling, Random cropping, Horizontal or vertical flipping, Adding noise (e.g., Gaussian noise), Changing brightness, contrast, or saturation

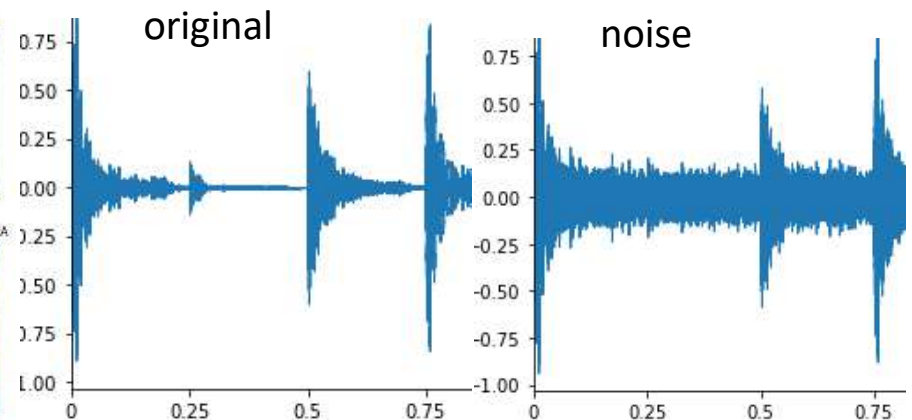
A  
P  
P  
l  
i  
c  
a  
t  
i  
o  
n  
s



(Self driving cars)



(Healthcare)



(Automatic Speech Recognition)

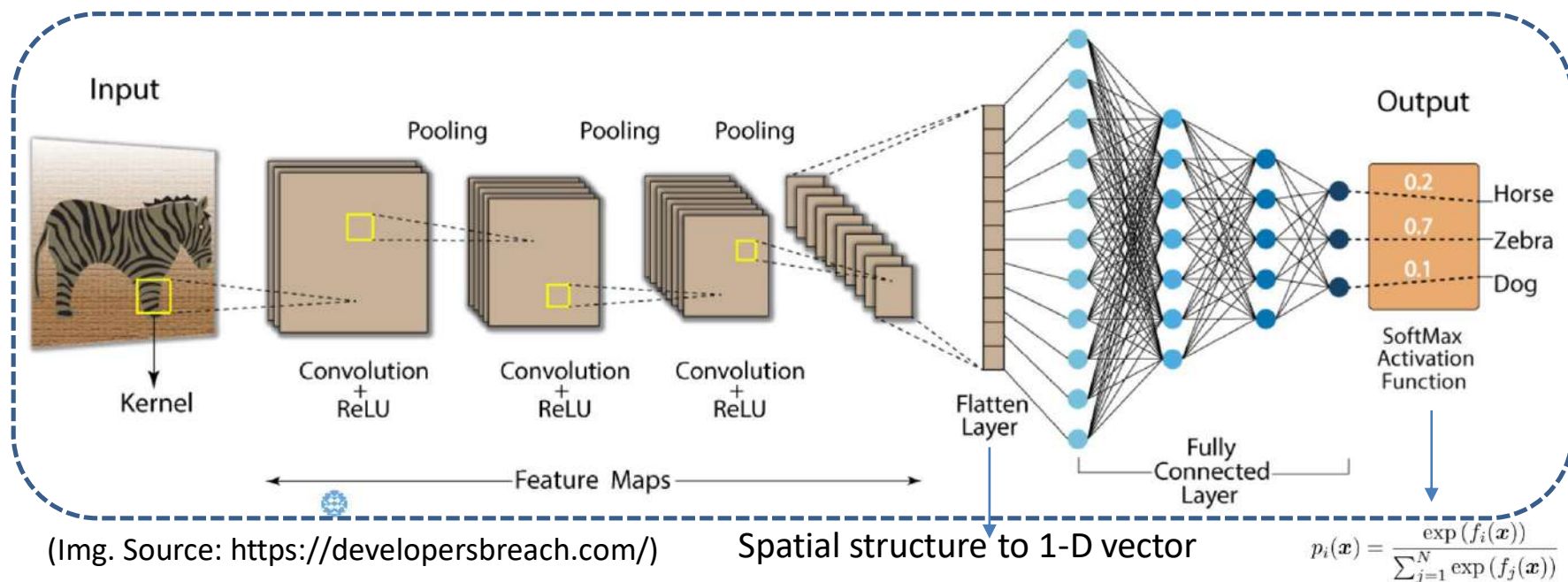
Img. Source: <https://www.datacamp.com/>

```
resize_and_rescale=keras.Sequential([ layers.Resizing(IMG_SIZE, IMG_SIZE), layers.Rescaling(1./255)])
```

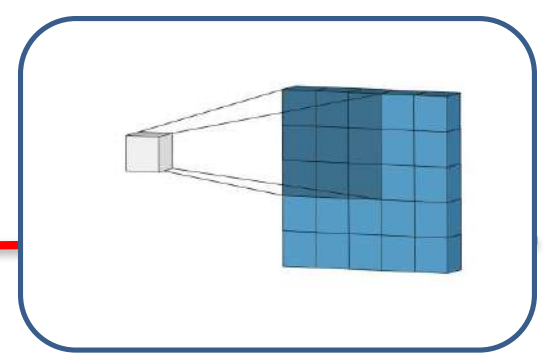


# Convolutional Neural Networks: Deep Learning

- So far...classified real values or discrete categories. What about Images & Sequences?
- Multi-layer Perceptrons (MLPs) are generally fully connected (each neuron in one layer is connected to every neuron in the subsequent layer).
- If Input: M units and Output N units, we need MXN connections. For an input image M of 256X256 = 65563 grayscale pixels, and output N of 1000 units, we would need 65 million connections.
- In Image data: We might want ‘Share structure property’ and ‘Invariance’ property to be encoded into the NN’s architecture. **CNNs to rescue.**

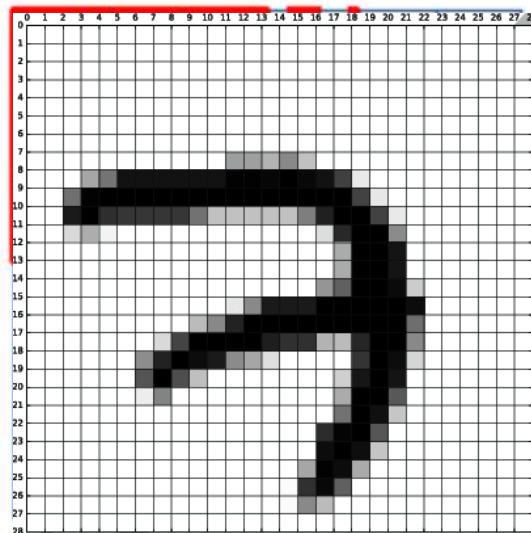


# CNN: Convolution Layer



- What is the need of Convolution Operation?
  - To extract features from input data by applying a **kernel** (also known as a **filter**) over the input.
  - When a **convolutional layer** is applied to an input image, the resulting feature maps often have **smaller spatial dimensions** compared to the input.

Dot product:  $(X * W) [i, j] = \sum_{m=-\infty}^{\infty} \sum_{n=-\infty}^{\infty} X[m, n] \cdot W[i - m, j - n]$



Original image (6X6)

Convolution Operation

Horizontal Sobel Kernel

$$\begin{bmatrix} -1 & 0 & 1 \\ -2 & 0 & 2 \\ -1 & 0 & 1 \end{bmatrix}$$

Kernel (3X3)

Vertical Sobel Kernel

$\times$  =

Corners or edges



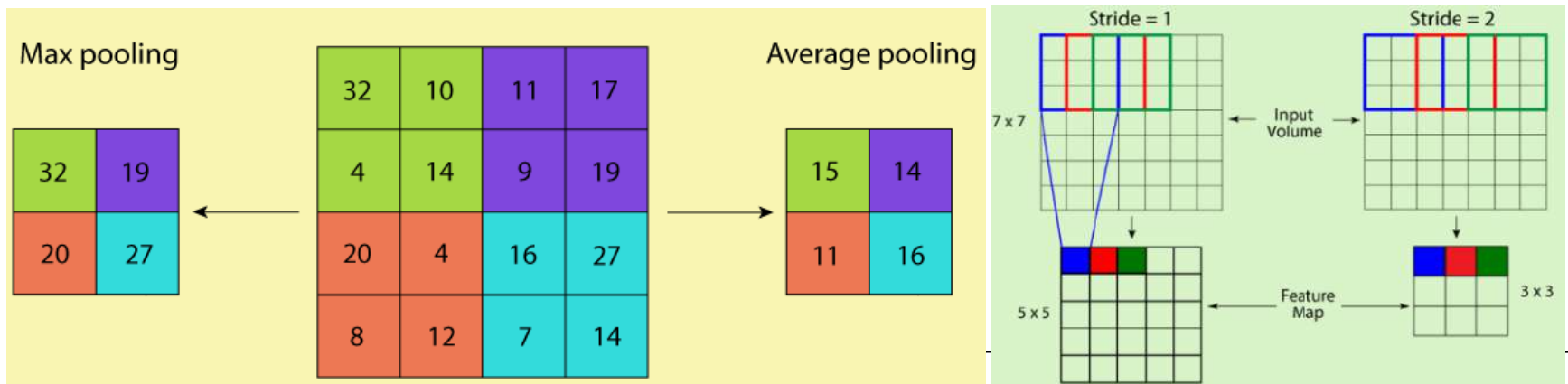
# CNN: Pooling (or Sub-sampling) Layer

**Spatial Invariance:** Pooling layers aggregate information from local neighborhoods of the input feature map, which helps in creating spatial invariance.

→ Spatial variations in the input (such as translation, rotation, or scaling) are tolerated to some extent, making the network more robust to variations in input data.

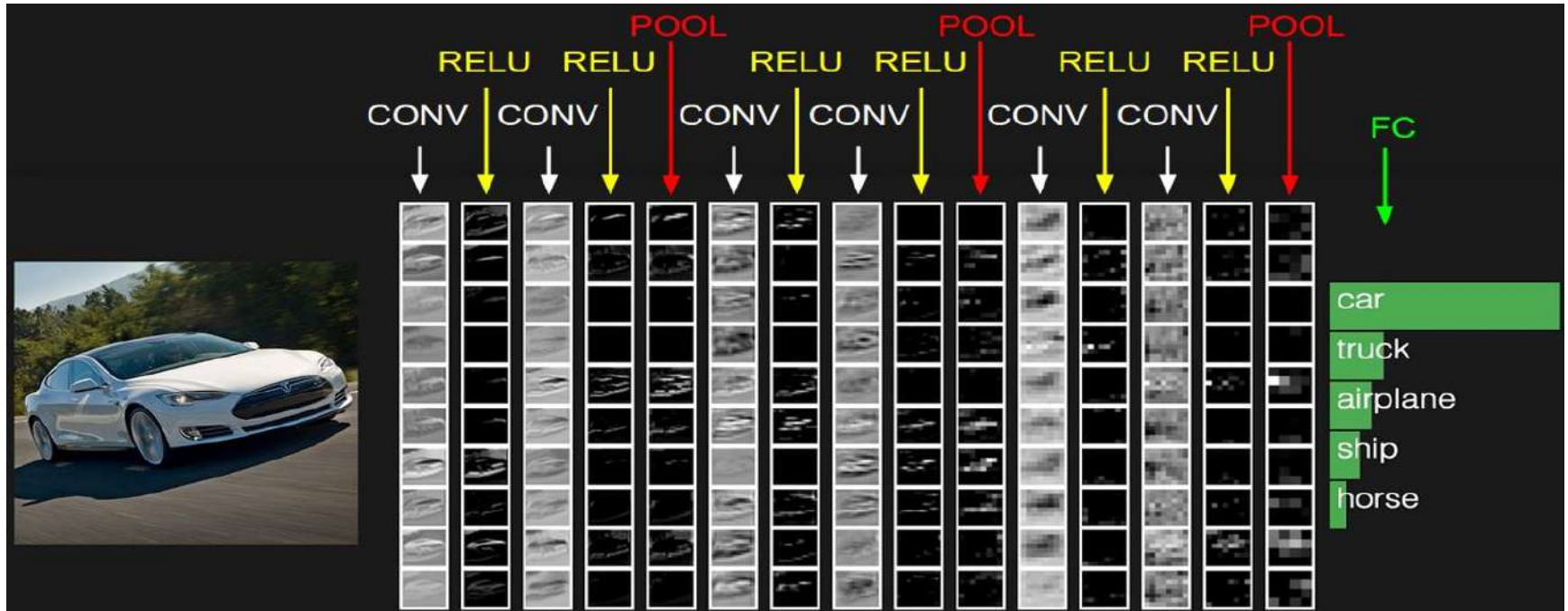
**Dimensionality Reduction:** downsamples the feature maps while retaining the important features. This reduces the number of parameters (weights) reducing the chances of **Overfitting**.

**Local feature detection:** Salient features within the local neighborhood is identified.



(Img. Source: <https://developersbreach.com/>)

# Classifying an Image: An example



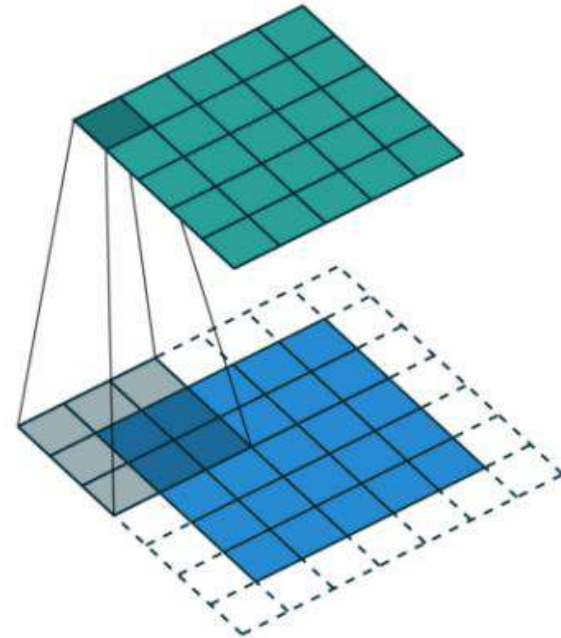
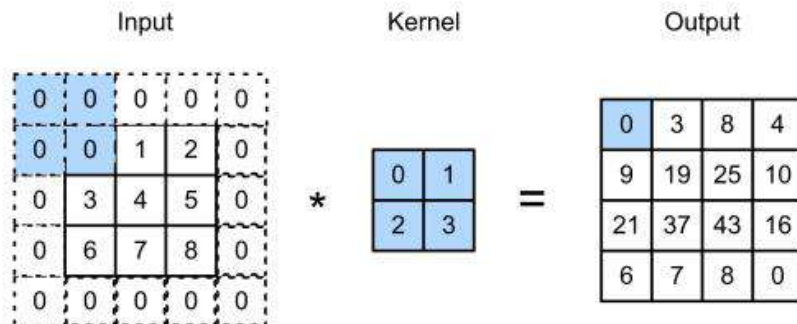
Acknowledgement: Md Mahin's presentation.



# CNN: Padding

---

- Padding in CNNs refers to the process of adding additional pixels around the input image **before applying** convolution operations.



- It controls the spatial dimensions of feature maps and prevents information loss at the borders of the image.

# Other types of Convolutional Networks: FCN, SSMD

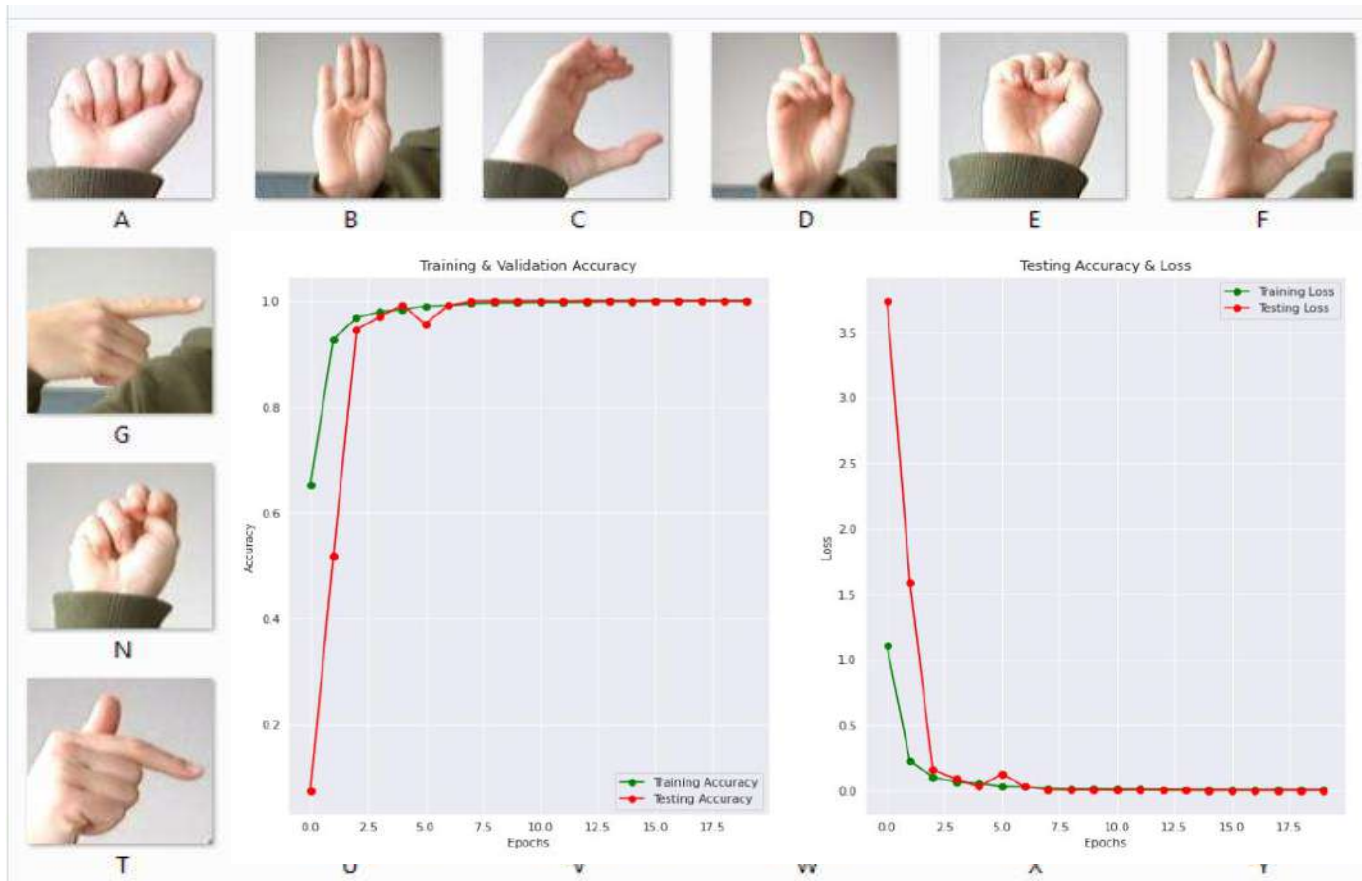
---

- **Semantic segmentation**, where the goal is to assign a class label to each pixel in the input image, the fully connected layer is not well-suited as they do NOT preserve spatial information.
- In **Fully Convolutional Networks** (FCNs), the final fully connected layer of the traditional CNN is replaced with **convolutional layers**. They allow the network to produce an output feature map with the same spatial dimensions as that of input.
- **Some Object Detection Networks**: In object detection architectures like YOLO (You Only Look Once) or SSD (Single Shot Multibox Detector), fully connected layers are often replaced by convolutional layers with spatial dimensions reduced to 1x1. Making network to predict **bounding boxes** and class probabilities at different spatial locations in the image.



# Mini-Project on hand gesture recognition using CNNs

Similar to LeNet



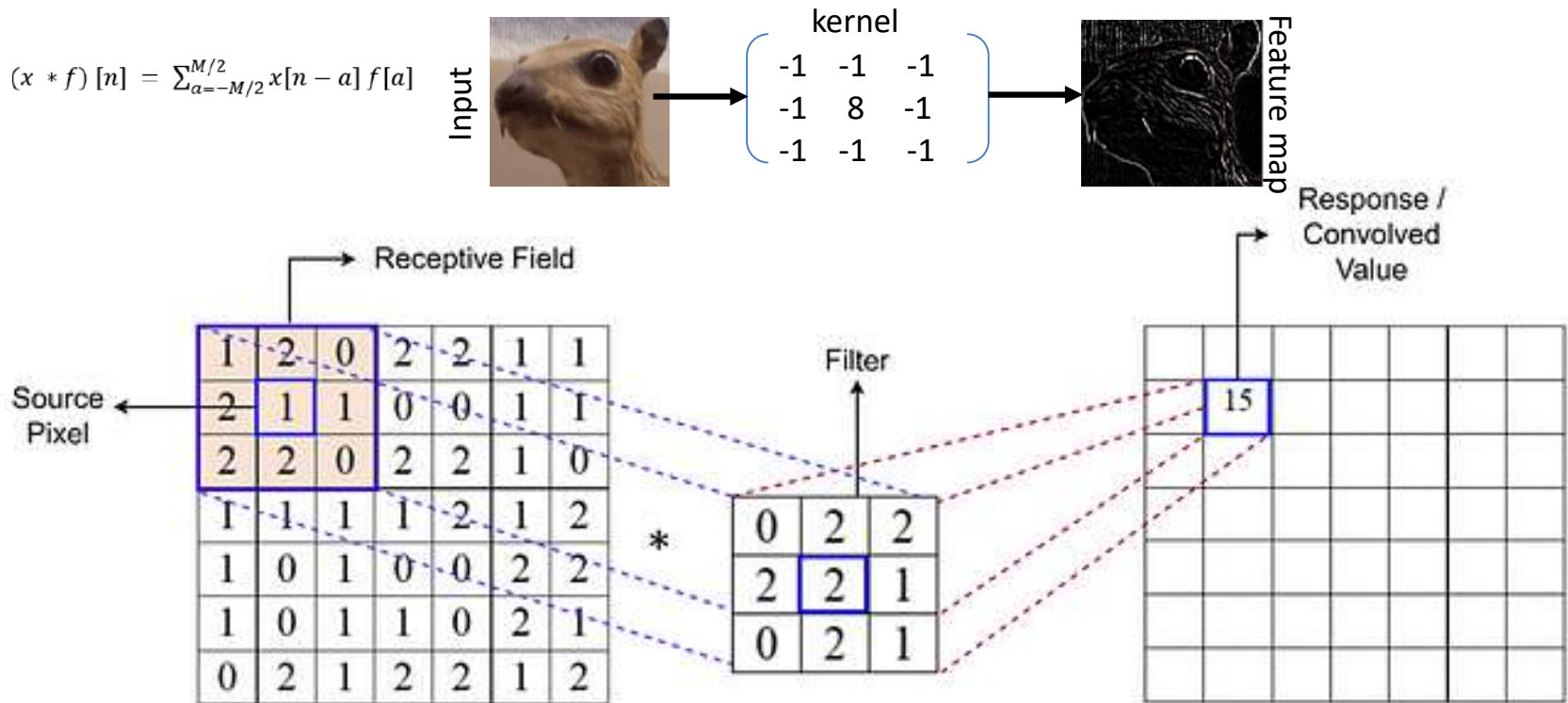
```
activation='relu',
```

```
activation='relu'),
```

Submission deadline: 30.04.2024

Many others: ResNet, AlexNet, ImageNet

# Recap: Convolution in CNNs

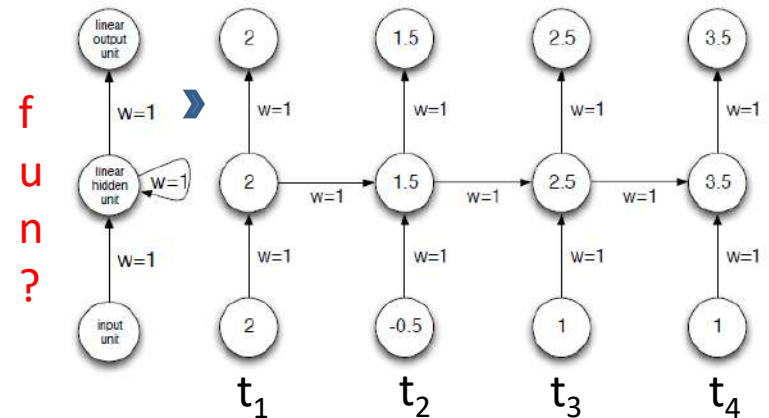
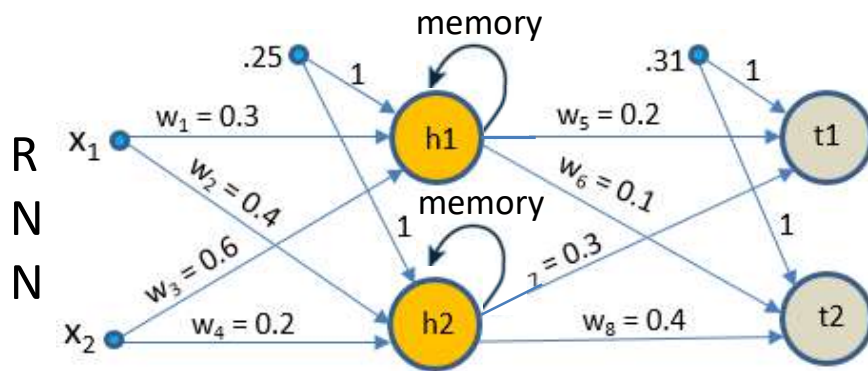


Convolution: A linear operation that computes the dot product of the local receptive field and the filter matrix.

It finds out the correction between the filter/ kernel with different parts of the image (when it slides over it), thereby learns important patterns or features.

# Recurrent Neural Networks (RNNs)

- Feed Forward Neural Networks are Acyclic where data passes from input to the output nodes and not vice versa.
  - Once the FFNN is trained, its state is fixed and does not alter as new data is presented to it. It **does not have memory**.



156.88 USD  
+0.88 (0.56%) ↑ today

Closed: 17 Apr, 7:58 pm GMT-4 • Disclaimer  
After hours 156.83 -0.050 (0.032%)

1D 5D 1M 6M YTD 1Y 5Y Max



Open	157.19	Mkt cap	1.94LCr	52-wk high	161.70
High	158.68	P/E ratio	28.91	52-wk low	103.27
Low	156.14	Div yield	-		

## Applications

39 °C | °F Precipitation: 20%  
Humidity: 31%  
Wind: 8 km/h

Temperature | Precipitation | Wind



1 pm	4 pm	7 pm	10 pm	1 am	4 am	7 am	10 am
Thu	Fri	Sat	Sun	Mon	Tue	Wed	Thu
39° 28°	39° 28°	38° 28°	38° 27°	38° 27°	39° 27°	39° 27°	39° 27°

Weather  
Thursday  
Partly Cloudy

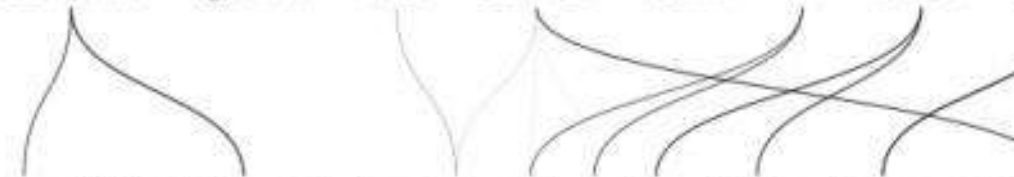


# More Applications...

---

## Machine Translation

Economic growth has slowed down in recent



Das Wirtschaftswachstum hat sich in den letzten Jahren verla  
Economic growth has slowed down in recent years



La croissance économique s' est ralentie ces dernières années

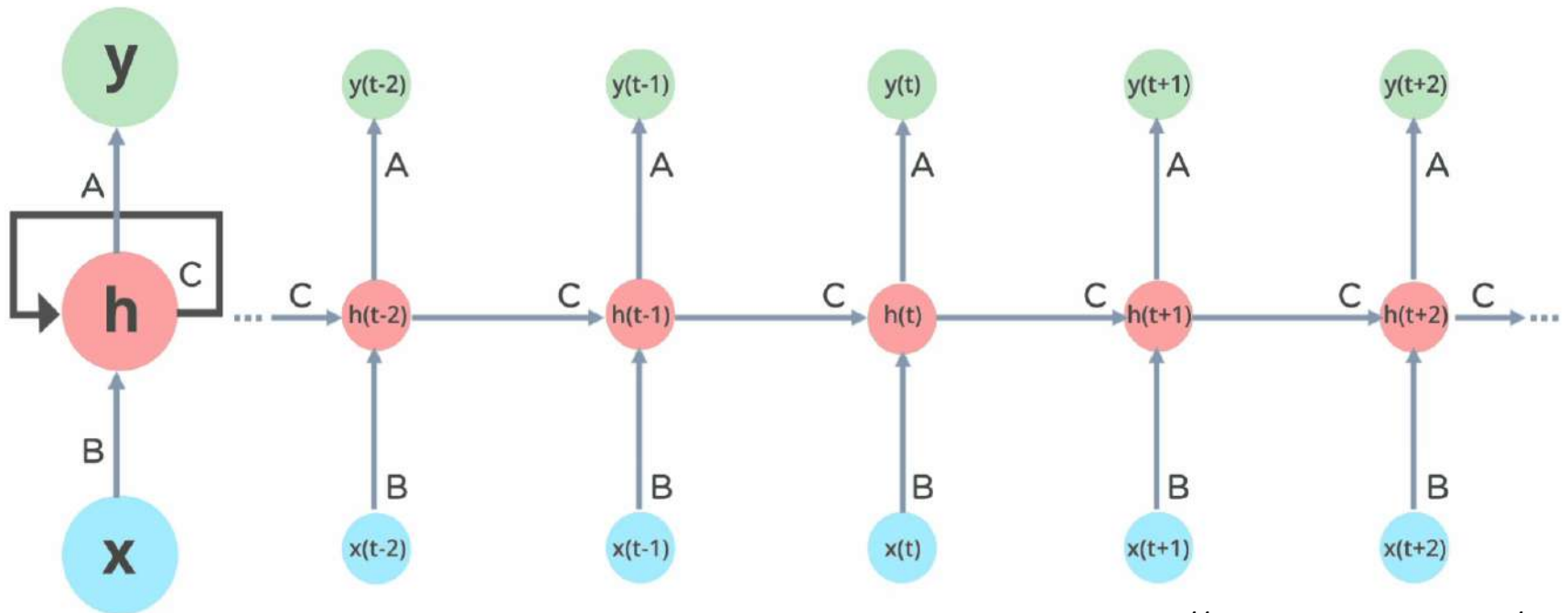


Video Classification

---

# Modelling sequential data: RNN

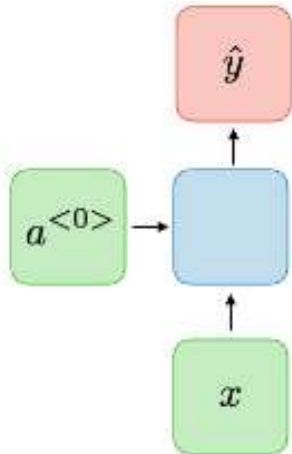
- A hidden state captures information about previous inputs. State is plugged back into itself.
- The hidden state is updated recurrently based on the current input and the previous hidden state.



<https://www.simplilearn.com/>

$g_1$  and  $g_2$ : Activation functions: Sigmoid or Tanh or ReLU

# RNN Architectures: One-to-One



One-to-one:  
Traditional  
NN

without  
considering  
temporal  
dependencies  
between  
reviews.

```
# Define the RNN model
model = tf.keras.Sequential([
    tf.keras.layers.Embedding(input_dim=len(tokenizer.word_index)+1, output_dim=16, ),
    tf.keras.layers.LSTM(32),
    tf.keras.layers.Dense(1, activation='sigmoid')
])

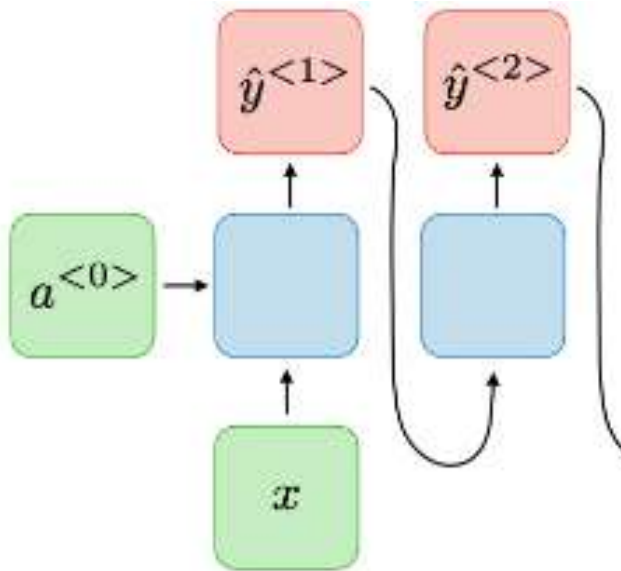
model.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy'])

# Train the model
model.fit(padded_sequences, labels, epochs=10, verbose=1)

# Inference
test_review = "This film was excellent!"
test_sequence = tokenizer.texts_to_sequences([test_review])
padded_test_sequence = pad_sequences(test_sequence, maxlen=max_length, padding='post')
prediction = model.predict(padded_test_sequence)

# Output the predicted sentiment
if prediction > 0.5:
    print("Positive sentiment")
else:
    print("Negative sentiment")
```

# RNN Architectures: One-to-Many



One-to-Many: Music generation

```
# Build the model
model = tf.keras.models.Sequential([
    tf.keras.layers.SimpleRNN(128, input_shape=(seq_length, num_chars), return_sequences=True),
    tf.keras.layers.TimeDistributed(tf.keras.layers.Dense(num_chars, activation='softmax'))
])

optimizer = tf.keras.optimizers.RMSprop(learning_rate=0.01)
model.compile(loss='categorical_crossentropy', optimizer=optimizer)

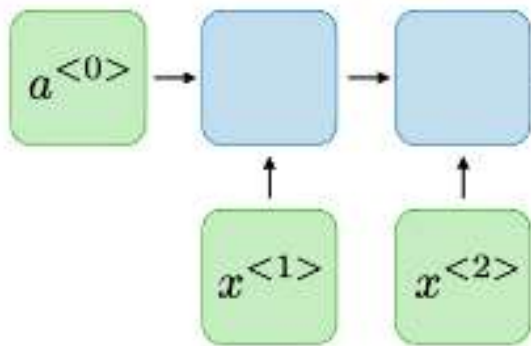
# Train the model
model.fit(X, y, batch_size=128, epochs=50)

# Generate text
def generate_text(seed_text, temperature=1.0):
    generated_text = seed_text
    for i in range(400):
        x_pred = np.zeros((1, seq_length, num_chars))
        for t, char in enumerate(seed_text):
            x_pred[0, t, char_to_idx[char]] = 1.0
        preds = model.predict(x_pred, verbose=0)[0][-1] # Take prediction from the last time step
        next_index = np.random.choice(len(chars), p=np.exp(preds) / temperature)
        next_char = idx_to_char[next_index]
        generated_text += next_char
        seed_text = seed_text[1:] + next_char
    return generated_text

# Generate text given an initial line
initial_line = "Tere sang jina yahan, tere sang mar jana"
generated_lyrics = generate_text(initial_line.lower())
print(generated_lyrics)
```

# RNN Architectures: Many-to-one

---



**Many-to-One:**  
Sentiment Classification

```
# Build the RNN model
model = tf.keras.models.Sequential([
    tf.keras.layers.Embedding(max_words, 16, input_length=max_sequence_length),
    tf.keras.layers.LSTM(32),
    tf.keras.layers.Dense(1, activation='sigmoid')
])

model.compile(optimizer='adam', loss='binary_crossentropy', metrics=['accuracy'])

# Train the model
model.fit(padded_sequences, labels, epochs=10, batch_size=32)

# Example text for prediction
example_text = "This product is fantastic!"

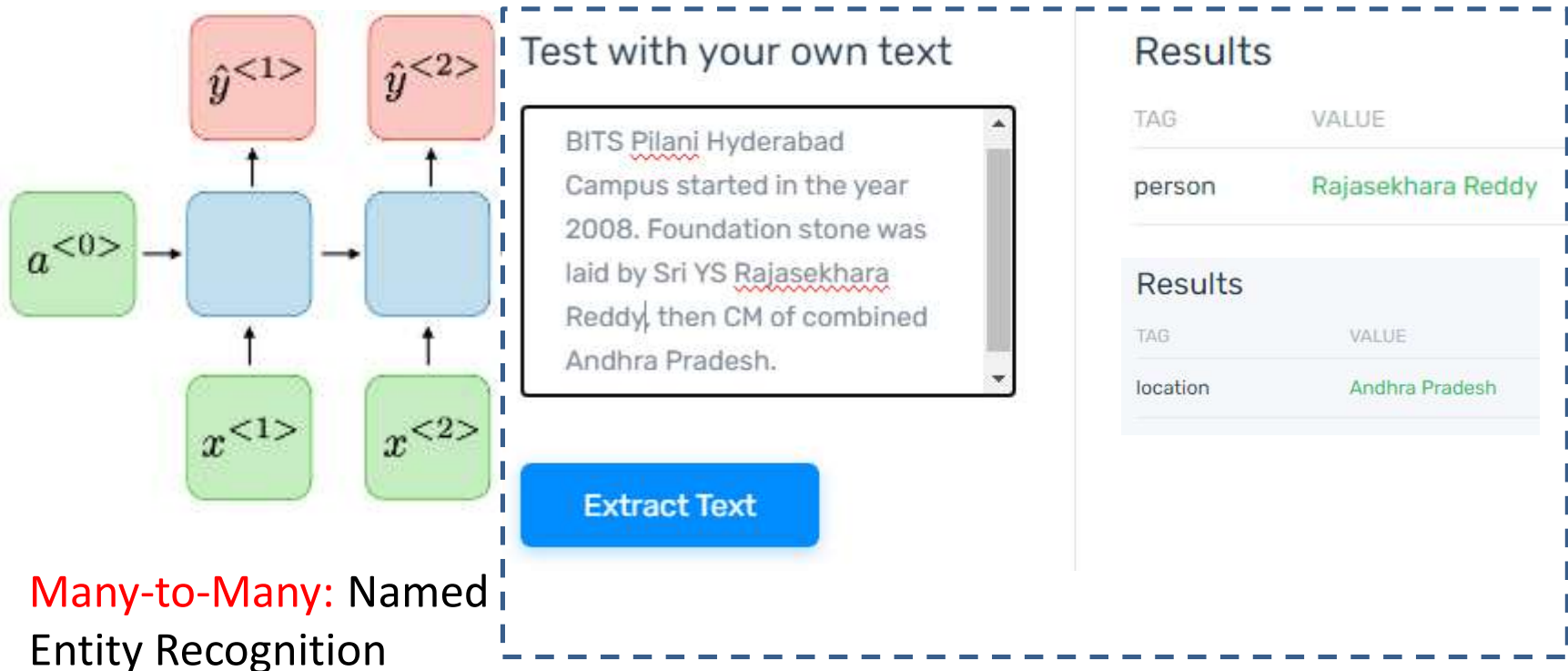
# Tokenize and pad the example text
example_sequence = tokenizer.texts_to_sequences([example_text])
padded_example_sequence = pad_sequences(example_sequence, maxlen=max_sequence_length)

# Predict sentiment
prediction = model.predict(padded_example_sequence)

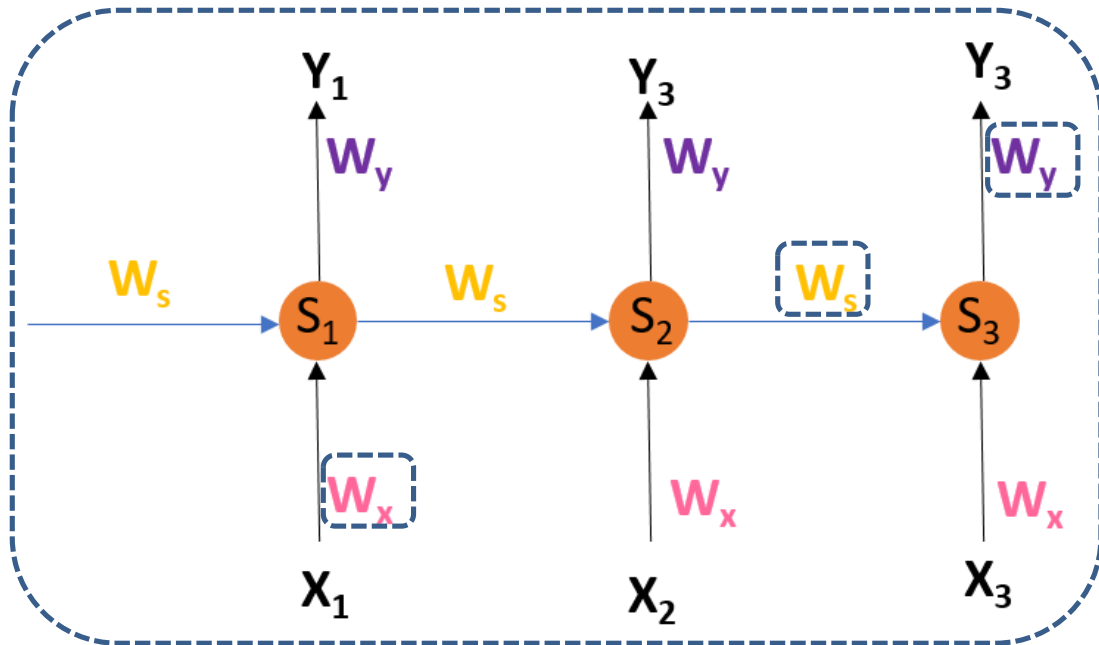
# Print prediction
if prediction[0] > 0.5:
    print("Positive Sentiment")
else:
    print("Negative Sentiment")
```



# RNN Architectures: Many-to-Many



# RNN Training: Backpropagation Through Time



$$E = (1/N) \sum_{i=1}^N (Y_i - \hat{Y}_i)^2$$

Adjusting  $W_y$  ( $E_3$  is a function of  $Y_3$  and  $Y_3$  is a function of  $W_y$ ):

$$\frac{\partial E_3}{\partial W_y} = \frac{\partial E_3}{\partial Y_3} \cdot \frac{\partial Y_3}{\partial W_y}$$

Similarly, now adjusting  $W_s$ :

$$\frac{\partial E_3}{\partial W_s} = \left\{ \frac{\partial E_3}{\partial Y_3} \cdot \frac{\partial Y_3}{\partial S_3} \cdot \frac{\partial S_3}{\partial W_s} \right\} + \left\{ \frac{\partial E_3}{\partial Y_3} \cdot \frac{\partial Y_3}{\partial S_3} \cdot \frac{\partial S_3}{\partial S_2} \cdot \frac{\partial S_2}{\partial W_s} \right\} + \left\{ \frac{\partial E_3}{\partial Y_3} \cdot \frac{\partial Y_3}{\partial S_3} \cdot \frac{\partial S_3}{\partial S_2} \cdot \frac{\partial S_2}{\partial S_1} \cdot \frac{\partial S_1}{\partial W_s} \right\}$$

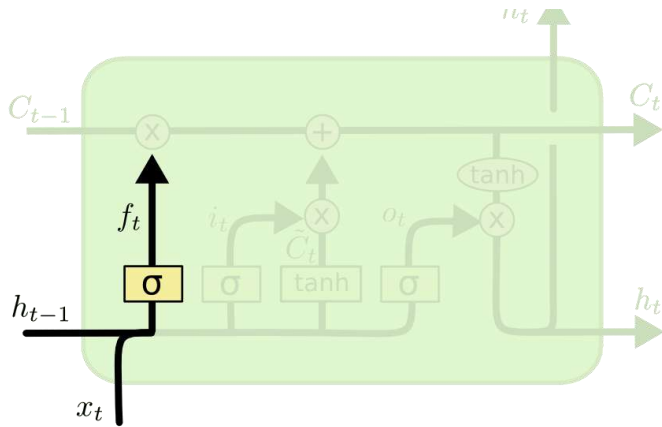
$E_3 \text{ --- } Y_3 \text{ --- } S_3 \text{ --- } W_s$ 
With respect to  $S_2$ 
With respect to  $S_1$

For  $W_x$ :

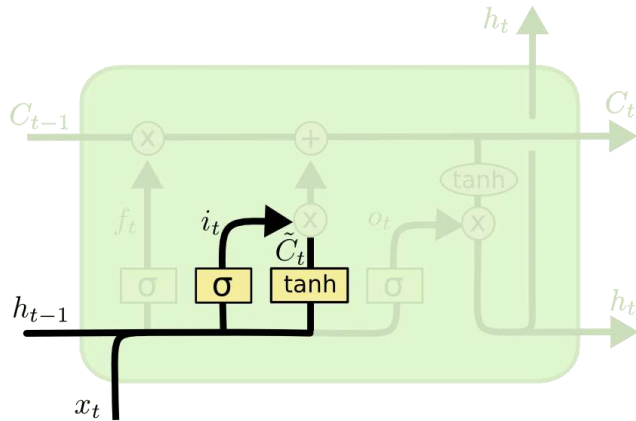
$$\frac{\partial E_3}{\partial W_x} = \left\{ \frac{\partial E_3}{\partial Y_3} \cdot \frac{\partial Y_3}{\partial S_3} \cdot \frac{\partial S_3}{\partial W_x} \right\} + \left\{ \frac{\partial E_3}{\partial Y_3} \cdot \frac{\partial Y_3}{\partial S_3} \cdot \frac{\partial S_3}{\partial S_2} \cdot \frac{\partial S_2}{\partial W_x} \right\} + \left\{ \frac{\partial E_3}{\partial Y_3} \cdot \frac{\partial Y_3}{\partial S_3} \cdot \frac{\partial S_3}{\partial S_2} \cdot \frac{\partial S_2}{\partial S_1} \cdot \frac{\partial S_1}{\partial W_x} \right\}$$

# LSTM: RNN

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



What info. is to be thrown away? → Forget gate (ex: gender of old subject)

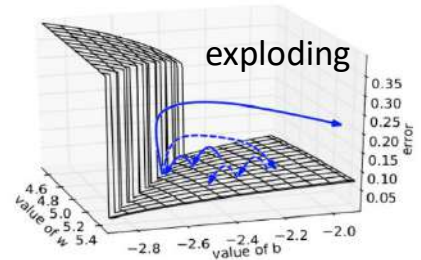
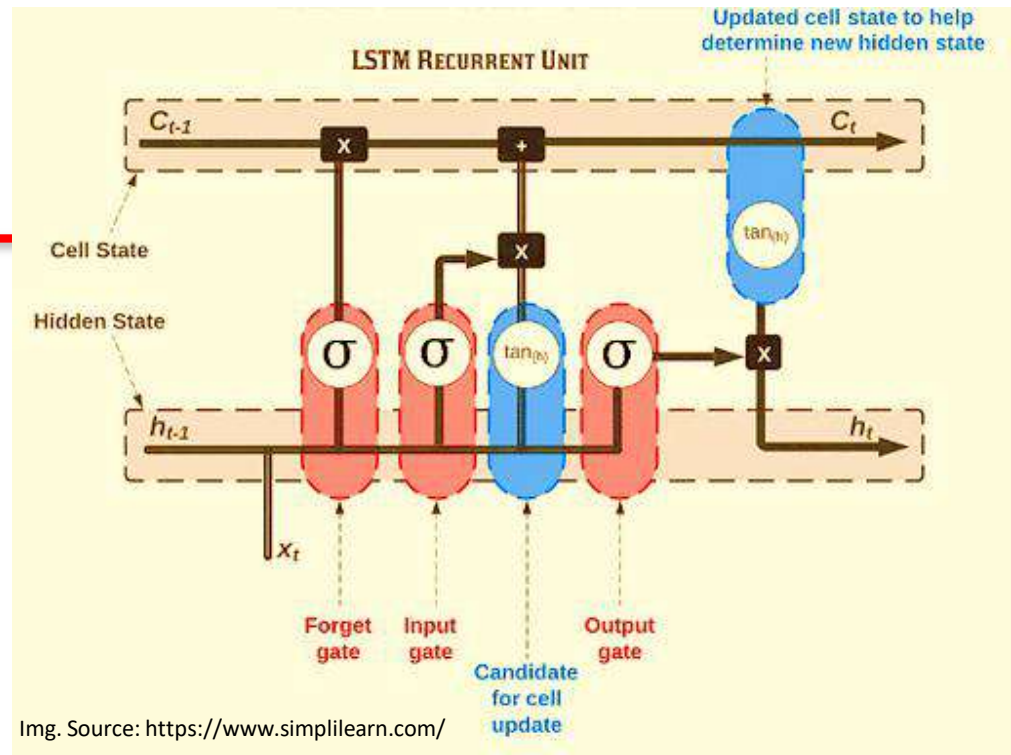


What new info. is to be stored?

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

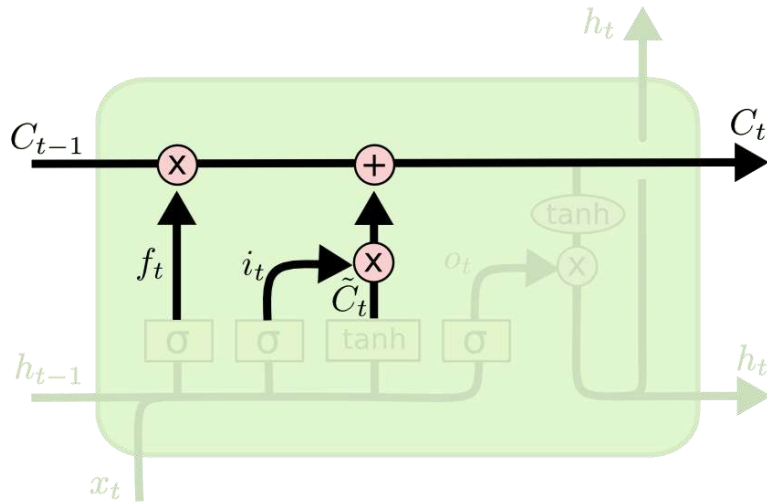
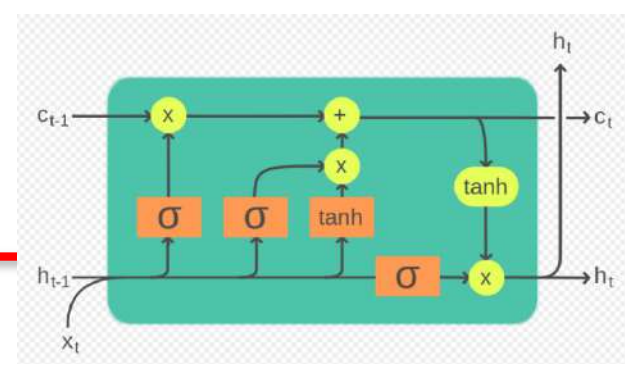
$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

May be gender of the new subject.



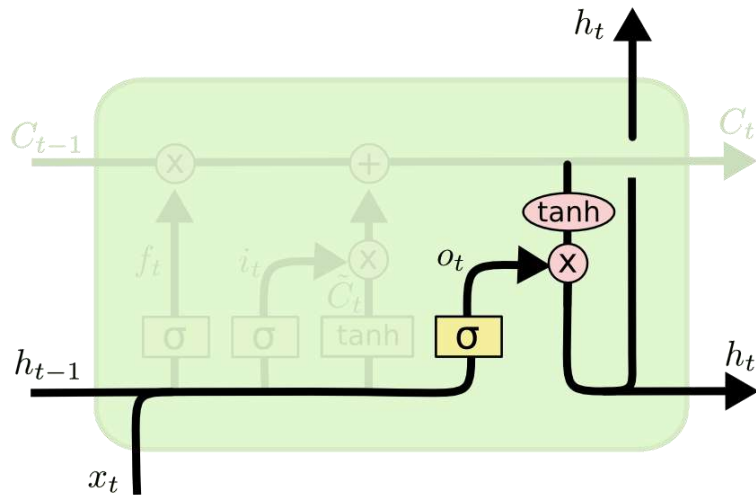
- **Problems with RNNs:** Vanishing and Exploding gradients. LSTMs solve vanishing prob. Exploding gradients prob. may be solved by Gradient clipping or regularization etc.

# Continued...



Update  $C_{t-1}$  into  $C_t$ .

$$C_t = \underbrace{f_t * C_{t-1}}_{\text{Forgetting the things}} + \underbrace{i_t * \tilde{C}_t}_{\text{New value}}$$



Output information relevant to a verb

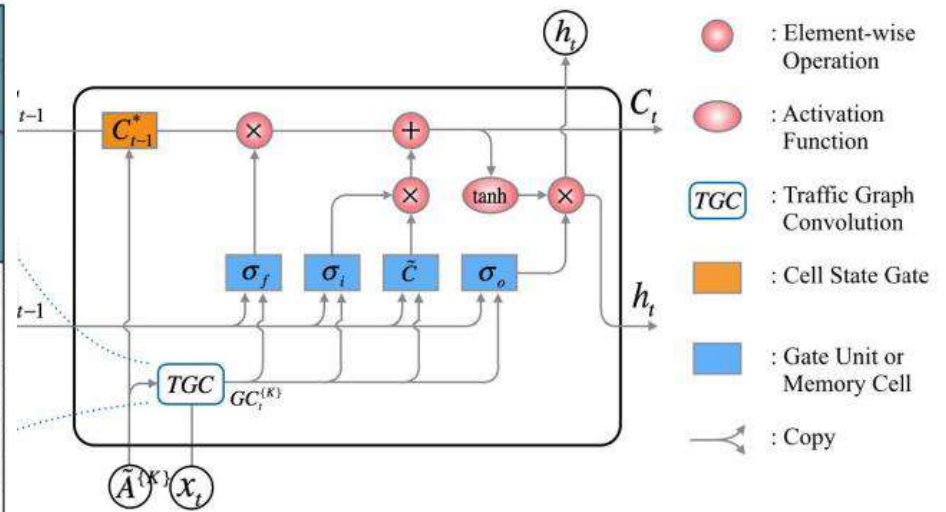
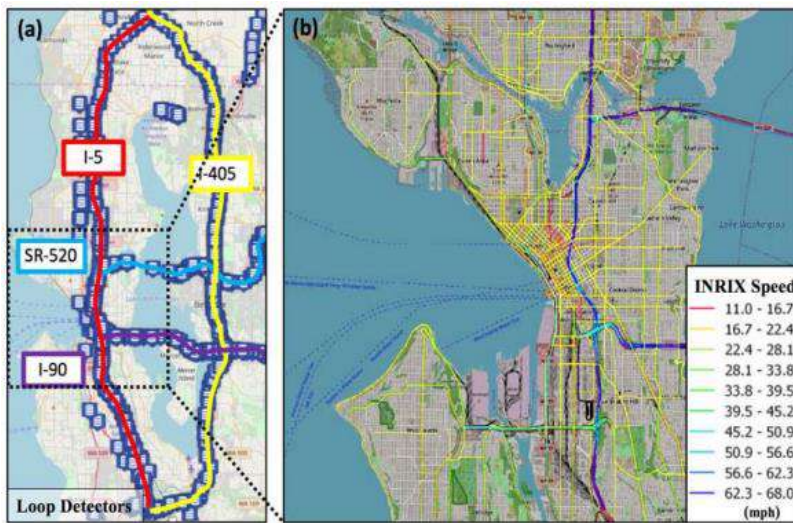
$$o_t = \sigma(W_o [h_{t-1}, x_t] + b_o)$$

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

# Traffic Graph Convolutional Recurrent Neural Network: A Deep Learning Framework for Network-Scale Traffic Learning and Forecasting

Zhiyong Cui<sup>1</sup>, Student Member, IEEE, Kristian Henrikson<sup>1</sup>, Ruimin Ke<sup>1</sup>, Student Member, IEEE, and Yin Hai Wang<sup>1</sup>, Senior Member, IEEE

**ConvLSTM:** Network wide traffic states are identified with most influential roadways.



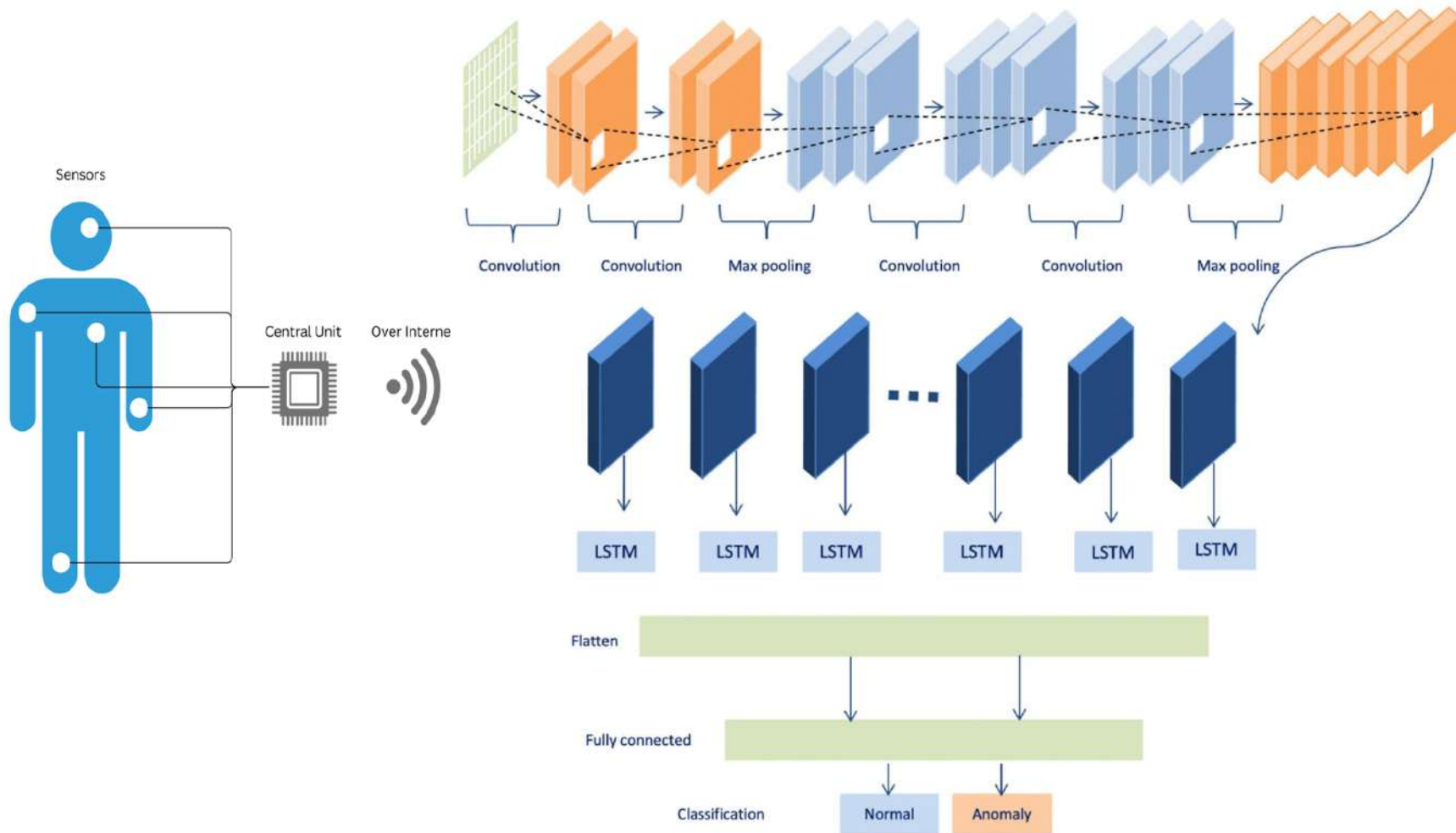
FFR: Free Flow Reachability (i.e vehicle speed) info,  $\tilde{A}$ : neighbourhood information,  $W$ : weights. Convolution layer is within the cell. Spatiotemporal dataset.



Article

# A Correlation-Based Anomaly Detection Model for Wireless Body Area Networks Using Convolutional Long Short-Term Memory Neural Network

Albatul Albattah <sup>1</sup> and Murad A. Kassam <sup>1,2,\*</sup>



Convolution + LSTM: genre of a movie by seeing the trailer (Horror or Detective)

# Autoregressive Models

- A generative model that generates new data points by regressing each observation on previous observations within the series.

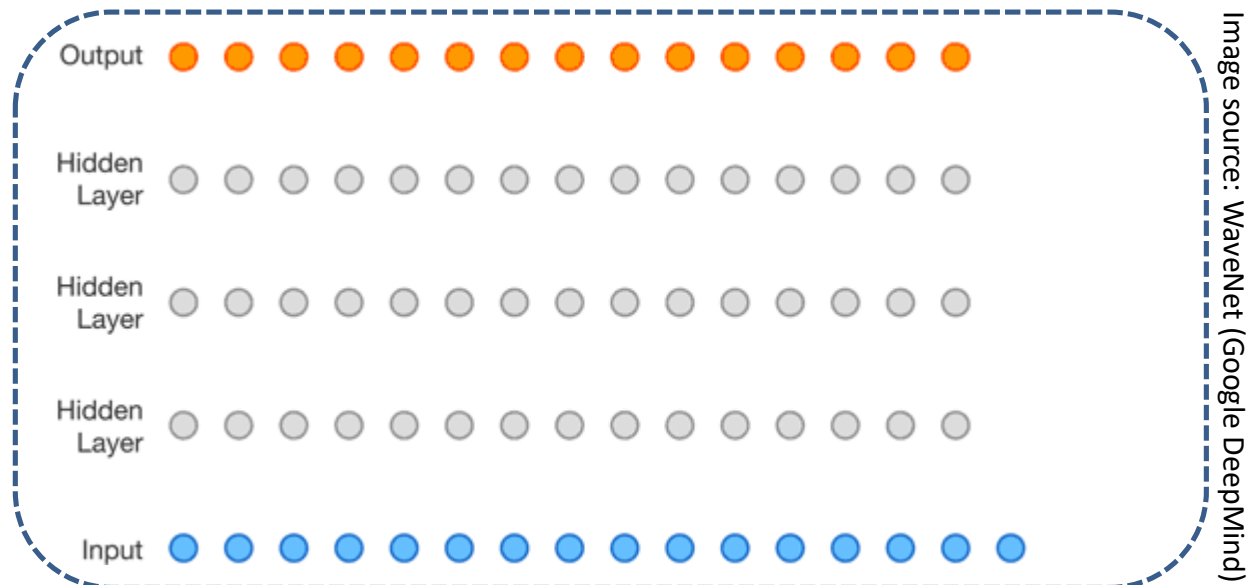
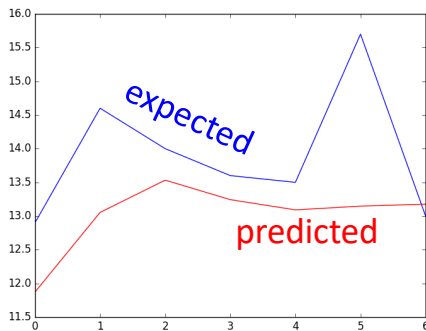
- Mathematically:  $Y_t = \Phi_0 + \Phi_1 \cdot Y_{t-1} + \xi_t$

$Y_t$  → Current value  
 $Y_{t-1}$  → Previous value

Where,  $\Phi_0, \Phi_1$  are coefficients to be estimated.  $\xi_t$  is error term at time 't'.

AR(1): Order 1

What would be AR(2), AR(3),...



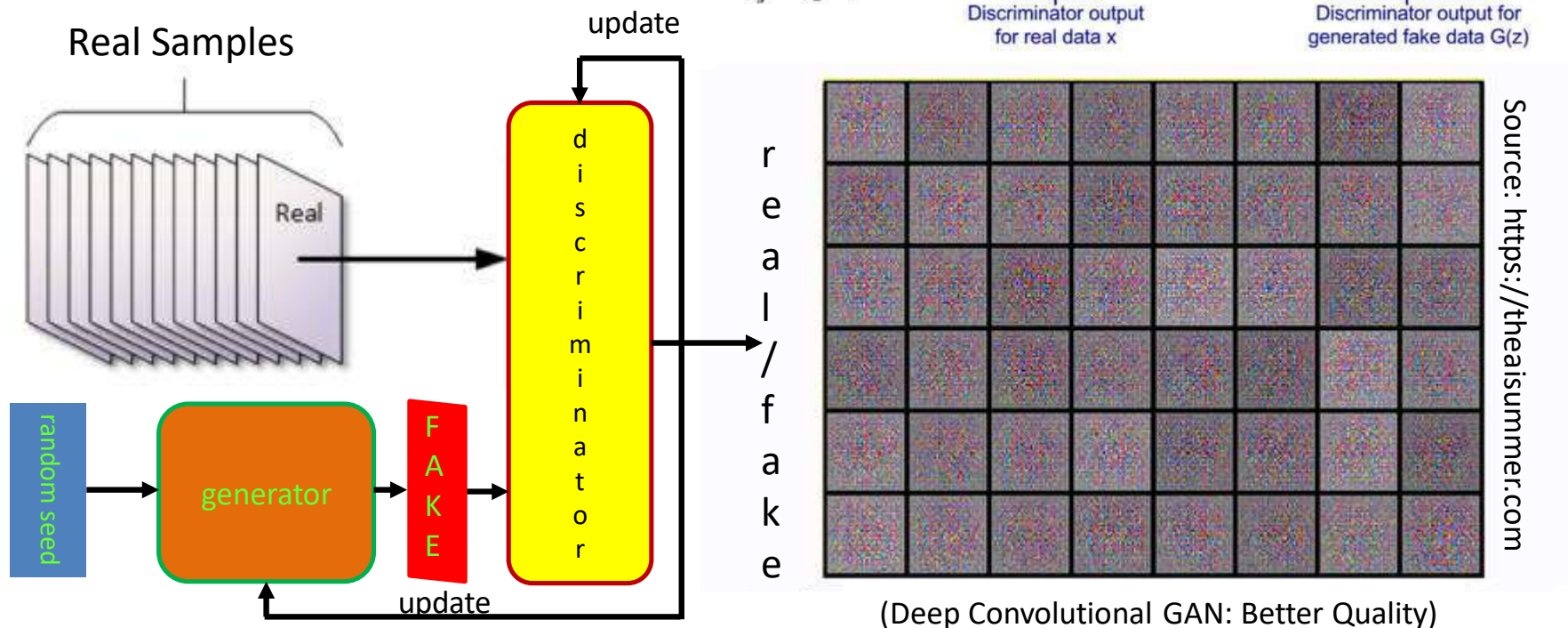
Probabilistically, Autoregression can be expressed as:  $p(x) = \prod_{i=1}^n p(x_i | x_1, x_2, \dots, x_{i-1}) = \prod_{i=1}^n p(x_i | x_{<i})$

# Generative Adversarial Network (GAN)

- GANs are neural networks (Deep ANNs) that learn to create synthetic data similar to some known input data. Or have the capability to generate new data.
- Ex: Admission office (Discriminator) checking the transcripts of newly admitted students( Generator)...

$$\min_{\theta_g} \max_{\theta_d} \left[ \mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

Discriminator output for real data x      Discriminator output for generated fake data G(z)



**Applications:** synthetic image/ video generation (deepfake), Image-to-image translation, Anomaly detection,...

---

Thank You!

---