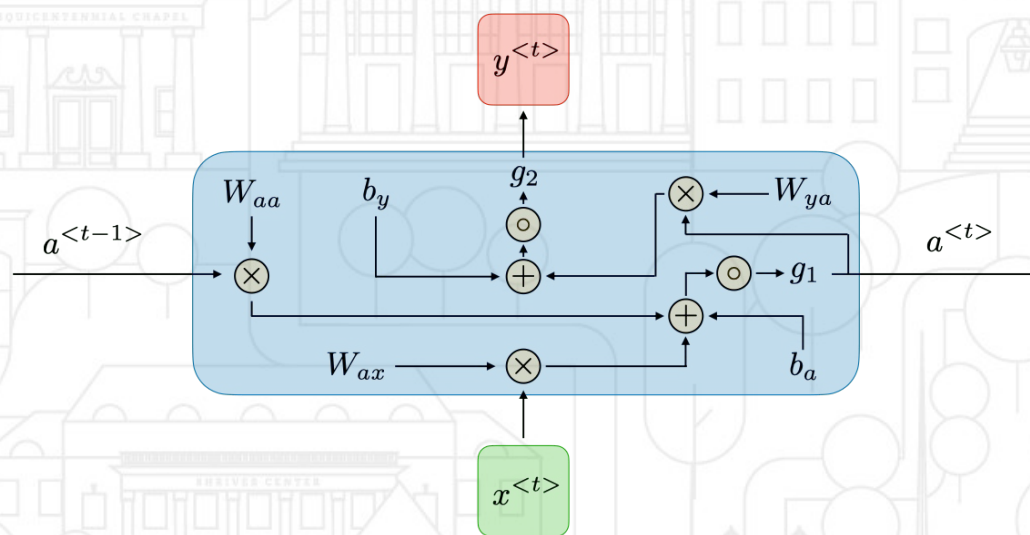# Module 6

Recurrent Neural Networks
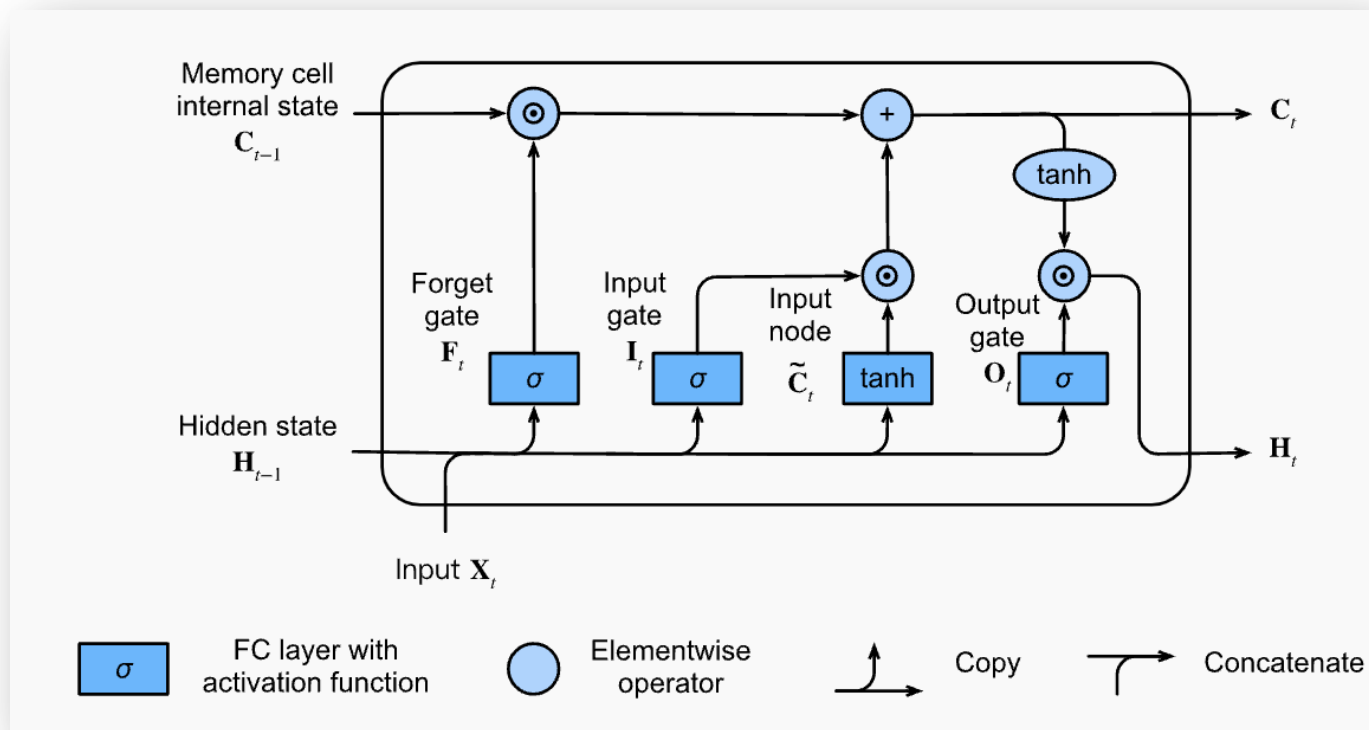
# DL: RNNs

# LSTM Gates

☑ Gates allow LSTMs to control short- and long-term memories



$$F_t = \sigma(X_t W_{xf} + H_{t-1} W_{hf} + b_f)$$

$$I_t = \sigma(X_t W_{xi} + H_{t-1} W_{hi} + b_i)$$

$$\tilde{C}_t = \tanh(X_t W_{xc} + H_{t-1} W_{hc} + b_c)$$

$$C_t = F_t \odot C_{t-1} + I_t \odot \tilde{C}_t$$

$$O_t = \sigma(X_t W_{xo} + H_{t-1} W_{ho} + b_o)$$

$$H_t = O_t \odot \tanh(C_t).$$

# Encoder-Decoder LSTM

❑ Seq2seq LSTM encoder-decoder networks
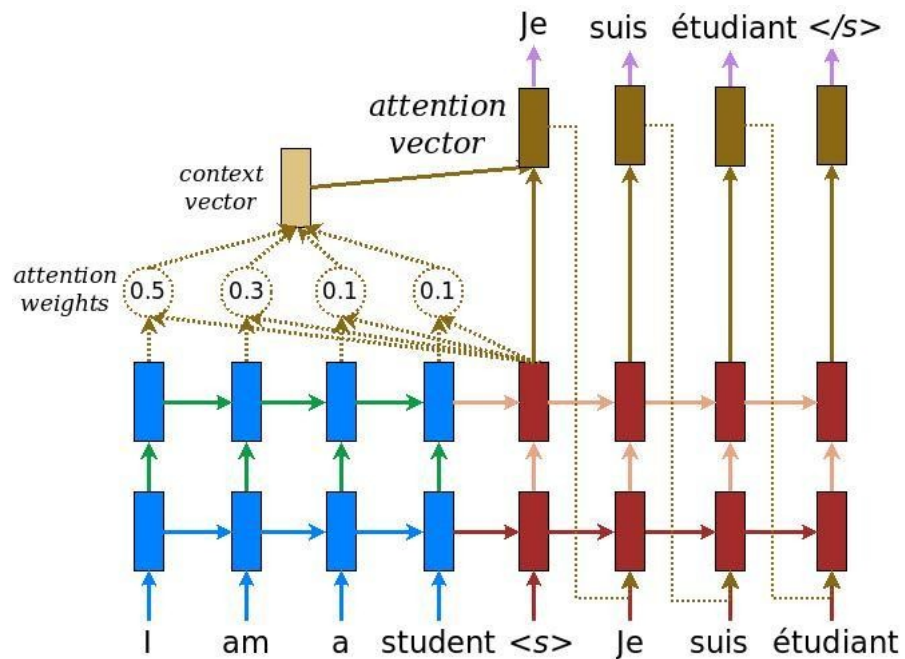
# Attention Layers

❏ Traditional models like RNNs and LSTMs treat all time steps equally.
❏ In reality, some moments matter more than others.
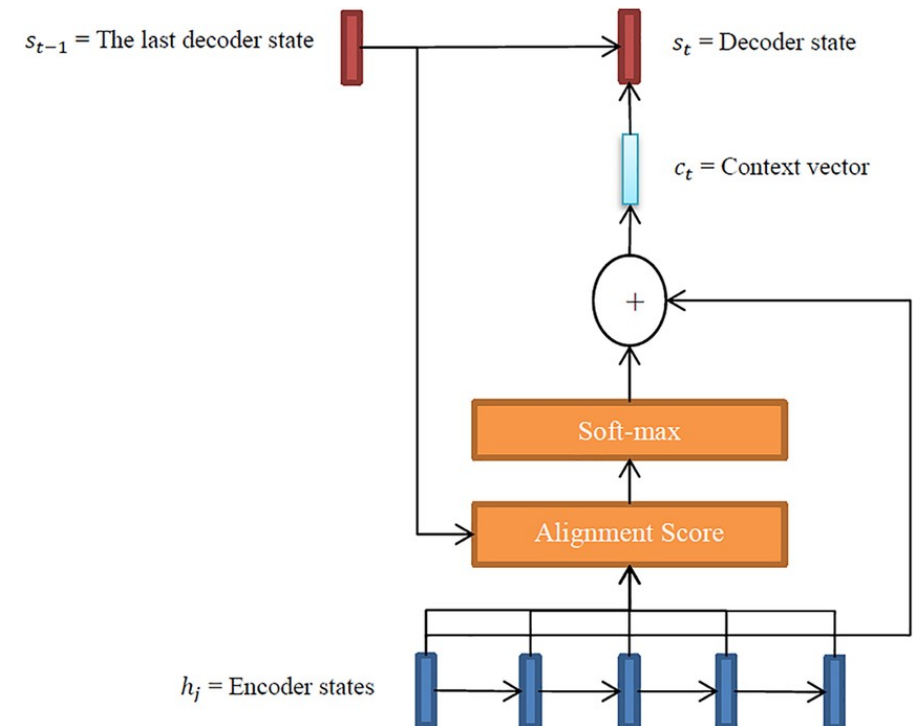❏ Goal: Learn to focus on important time steps dynamically.

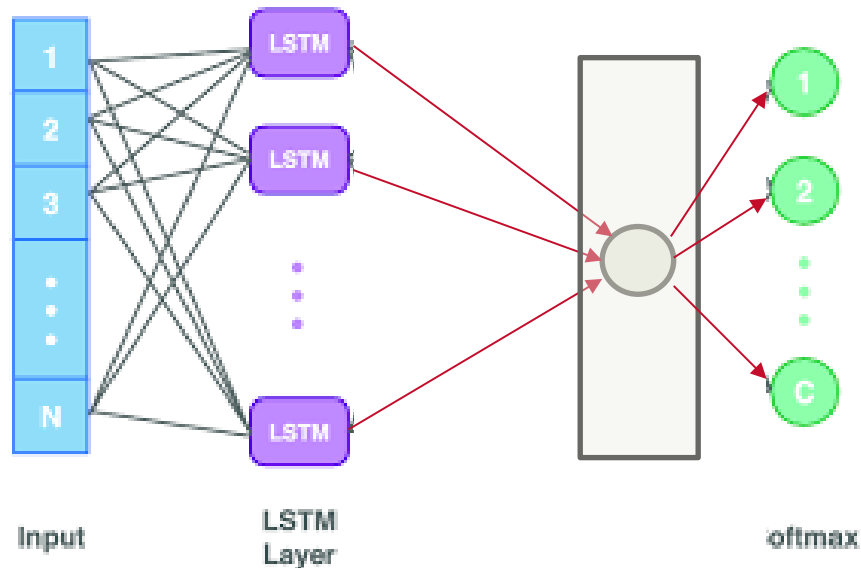# What is Attention?

❑ A mechanism that lets models learn where to focus.
❑ Assigns a weight to each input element.
❑ Outputs a context vector summarizing important parts.

# How Attention Works

❑ Score each time step with a small neural network.
❑ Apply softmax to get attention weights.
❑ Multiply each hidden state by its weight.
❑ Sum to get the context vector.

# Attention Layers (Business Uses)

- ❑ Besides the regular use of LLMs
    - ❑ Forecasting
    - ❑ Past sales, promotions, holidays, weather
    - ❑ Attention focuses on:
        - ❑ Holiday seasons
        - ❑ Promotional weeks
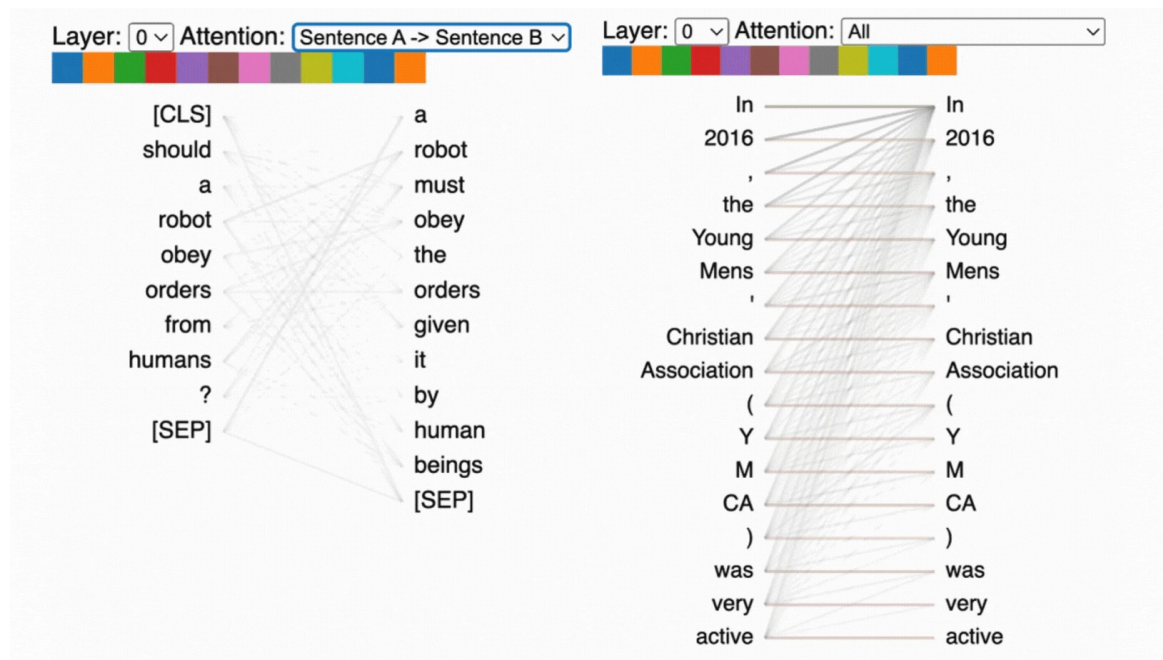        - ❑ Special weather events

# Attention Layers Anatomy

❑ Example of Attention Layers:
   ❑ **Dense(1):** Assigns a score to each time step.
   ❑ **Softmax:** Normalizes into probabilities.
   ❑ **Reduce_sum:** Dot Product creates the context vector.

```python
attention_scores = Dense(1, activation='tanh')(lstm_out)
attention_weights = Activation('softmax')(attention_scores)
context_vector = tf.reduce_sum(attention_weights * lstm_out,
axis=1)
```

# Benefits of Attention Layers

- ❏ Improved Accuracy: Focus on key time steps.
- ❏ Interpretability: Understand what the model "looks at."
- ❏ Scalability: Easier learning for long sequences.

# Python

# Multi-Head Attention Layers

- ❑ Introduced in Transformer models.
- ❑ Applies multiple attention heads in parallel.
- ❑ Each head learns to focus on different aspects of the input.
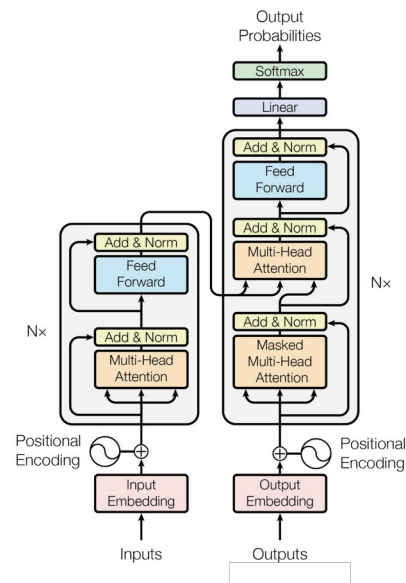- ❑ Outputs from all heads are concatenated and linearly transformed.

| Single Head | Multi-Head |
|---|---|
| One perspective on sequence | Many parallel perspectives |
| May miss complex patterns | Captures diverse relationships |
| Limited capacity | Higher representational power |

# Multi-Head Attention Layers

❑ Introduced in Transformer models.
❑ Applies multiple attention heads in parallel.
❑ Each head learns to focus on different aspects of the input.
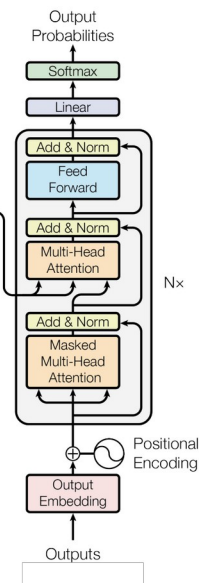❑ Outputs from all heads are concatenated and linearly transformed.

**BERT**

Encoder

**GPT**

Decoder

Python