

# RNN (Recurrent Neural Networks) [used in NLP]

- is a type of Neural Network designed to process sequences of data (like text, speech, or time series) by incorporating an internal "memory" of previous inputs. Unlike traditional neural networks, which treat inputs independently, the RNN uses information from prior steps to influence the current output
- How they work: The core mechanism that gives RNN's their memory is the hidden state and Feedback loops
  - 1) sequential processing: RNN processes data one element (word, fractal, etc) at a time, in order. The output at a specific time step depends on both the current input and information retained from the previous time steps
  - 2) The hidden state (memory) at each step, the network updates an internal hidden state vector. This hidden state acts as a summary or "memory" of all the sequence data the network has processed so far
  - 3) Feedback loop and weight sharing: The crucial part is that this updated hidden state is fed back into the network as an input for the next time step. The same set of weights and biases are used (shared) across all time steps, allowing the network to learn generalized patterns across sequences regardless of their position!
  - 4) Contextual understanding: This looping mechanism allows the network to build a contextual understanding of the entire sentence. For example when predicting the next word in a sentence, the RNN uses the history of the previous words to make an informed prediction, rather than just using the single current word.

Example: To predict the word "red" in sequence "Apple is ..."

- Timestep 1: when the RNN processes "Apple", it stores a representation in its memory (hidden state) a numerical vector
- Timestep 2: when the RNN processes "is", it recalls "Apple" from its memory and understands the full context "Apple is". To better predict that "red" is likely next word
- Note: storing of "Apple" happens when hidden state vector is updated by combining the input vector for "Apple" with previous hidden state using weighted matrix + activation function (non-linear)
- Type 3: to solve RNN issues like not remembering information in long sequences (gradient vanishing)
- LSTM (long short term memory): uses a system of gates (input, forget and output) to control which info in the memory can be kept or discarded allowing to capture long term dependencies effectively
- GRU (Gated Recurrent Unit): A simpler more efficient alternative to LSTM with only 2 gates (reset and update) that also solves the long term dependency problem. An update gate controls (reset and update) that also solves the long term dependency problem. An update gate decides how much of the input should be combined with existing memory cell, reset gate decides how much of past memory to forget before considering new input.