

# Recurrent Neural Networks

**Tianxiang (Adam) Gao**

February 22, 2024

# Outline

- 1 Speech and Language Problems
- 2 Recurrent Neural Networks (RNNs)
- 3 Stabilize RNNs Learning
- 4 Word Embedding

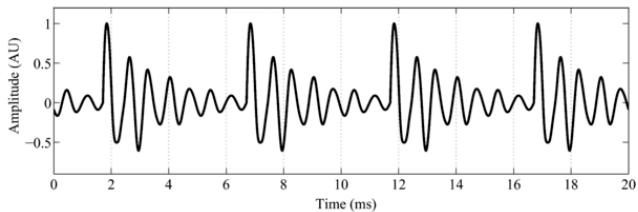
## Recap: Learning with CNNs

- **$1 \times 1$  Convolution:** Learns high-level patterns by combining multiple basic patterns.
- **Classic CNNs:** Inception with various setups, and MobileNet using depthwise separable convolution.
- **Transfer Learning:** Fine-tune a large, pretrained model on a smaller dataset using a lower learning rate to learn task-specific features.
- **Data Augmentation:** Increases dataset diversity and reduces overfitting (e.g., flips, random cropping, color adjustments, mixups).
- **Object Detection:** Uses  $1 \times 1$  convolution to implement fully connected layers (FC). The YOLO algorithm learns both class distribution and bounding boxes by leveraging object localization.
- **Semantic Segmentation:** Assigns a class label to each pixel. UNet combines lower- and higher-level features using an encoder-decoder architecture.
- **Face Recognition:** Learns a similarity function explicitly (via triplet loss) or implicitly (via siamese networks).
- **Neural Style Transfer:** Minimizes content and style loss simultaneously, where style is defined as the correlation between different channels.

# Outline

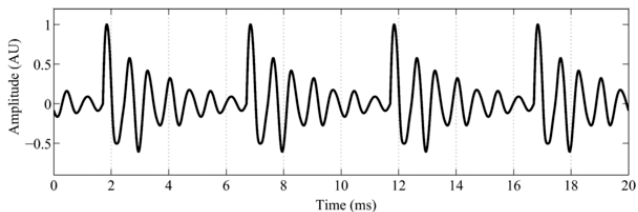
- 1 Speech and Language Problems
- 2 Recurrent Neural Networks (RNNs)
- 3 Stabilize RNNs Learning
- 4 Word Embedding

# Speech Emotion Recognition



- **Input:** The raw audio waveform
- **Output:** Angry? Multi-class classification

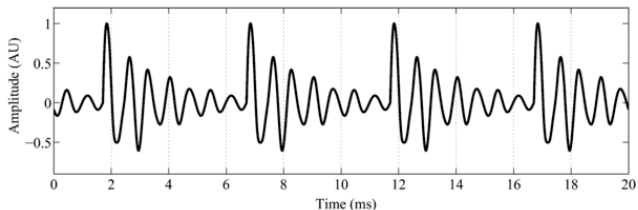
# Speaker Identification and Verification



- **Input:** The raw audio waveform
- **Output:** Classification or similarity learning



# Speech Recognition



- **Input:** The raw audio waveform
- **Output:** A **sequence** of text or words
- **Model:** It is framed as a sequence-to-sequence problem.

*“Good morning, everyone! Today, we’ ll be discussing sequential data and sequential models and their use in tasks like speech recognition and language processing. Feel free to ask questions during the lecture.”*

# Sentiment Analysis

*"I recently got a new smartphone, and I'm thrilled! The camera is amazing, and the battery lasts all day. I'm a bit disappointed with the fingerprint sensor, but overall, it's a great phone, and I'm happy with my purchase."*



Negative



Neutral



Positive

- **Input:** A sequence of text.
- **Output:** sentiment classification, e.g., tone of the text.



# Machine Translation

*"I recently got a new smartphone, and I' m thrilled! The camera is amazing, and the battery lasts all day. I' m a bit disappointed with the fingerprint sensor, but overall, it' s a great phone, and I' m happy with my purchase."*

“我最近买了一部新智能手机，感到非常兴奋！相机效果非常棒，电池可以用一整天。不过，指纹传感器有点让我失望，但总体来说，这是一部很不错的手机，我对这次的购买非常满意。”

- **Input:** A sequence of text in the source language
- **Output:** A sequence of equivalent text translated into the **target language**
- **Model:** seq2seq task and language model

# Chatbots

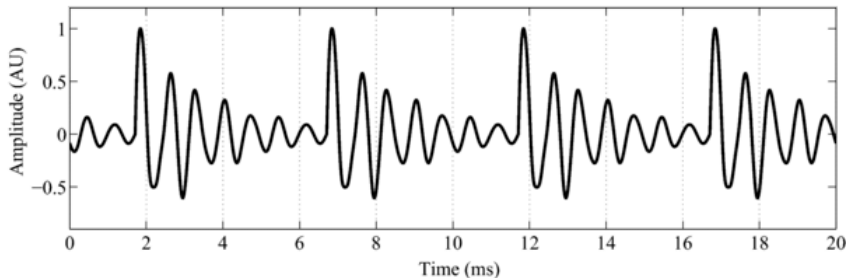
who are you?



I am Vita, your teaching assistant for CSC 578, a deep learning course. My role is to assist with questions about the course, covering topics like neural networks, optimization, CNNs, and more, as well as logistical queries about the class. If you have any course-related questions, feel free to ask!

- **Input:** A sequence of text or tokens
- **Output:** A sequence of response text
- **Presentation:** Word embeddings

# The Raw Audio Waveform



- **Define:** A 1D signal array represents the amplitude of the sound over time
- **Sampling Rate:** The number of samples per second (e.g., 16 kHz or 44.1 kHz).
- **Amplitude:** The “loudness” of the sound, represented either by integers ranging from -32,768 to 32,767 (for 16-bit audio) or normalized between -1.0 and 1.0 for floating-point.
- **Example:** A 10-second audio waveform sampled at 16kHz would result in a 1D array with  $16,000 \times 10 = 160,000$  amplitude values, each representing the sound amplitude at a specific point in time.

# One-Hot Encoding

- **Define:** Each word in a **vocabulary** is represented by a **binary** one-hot vector.
- **Sequence Data:**

$x^{(1)}$	"The dog chases the cat."
$x^{(2)}$	"A bird flies over the dog."
$x_3$	"The mouse hides from the cat."
$\vdots$	$\vdots$
$x^{(1000)}$	"The cat watches the fish swim."

- **Vocabulary:**

Word	Word Index	One-Hot Encoding
a	1	[1, 0, 0, ..., 0]
bird	2	[0, 1, 0, ..., 0]
cat	3	[0, 0, 1, ..., 0]
$\vdots$	$\vdots$	$\vdots$
zoo	10,000	[0, 0, 0, ..., 1]

- **Special words:** start of a sentence (**SOS**) and end of a sentence (**EOS**) with extra indices.

# Outline

- 1 Speech and Language Problems
- 2 Recurrent Neural Networks (RNNs)**
- 3 Stabilize RNNs Learning
- 4 Word Embedding

## Challenges in Text Data

- **Input Sequence:** Each input sequence is represented as

$$\mathbf{x}^{(i)} = \{\mathbf{x}_1^{(i)}, \dots, \mathbf{x}_{T_i}^{(i)}\} = \{\mathbf{x}_t^{(i)}\}_{t=1}^{T_i}, \quad \forall i \in [N]$$

where  $\mathbf{x}_t^{(i)}$  is the input at time step  $t$  of the  $i$ -th sequence,  $T_i$  is the sequence length, and  $N$  is the total number of sequences.

- **Training Dataset:**

$$\mathcal{D} = \{\mathbf{x}^{(i)}, \mathbf{y}^{(i)}\}_{i=1}^N$$

where  $\mathbf{y}^{(i)}$  is the target sequence corresponding to the input  $\mathbf{x}^{(i)}$ .

- Using one-hot encoding, each  $\mathbf{x}_t = \mathbf{e}_i \in \mathbb{R}^V$ , where  $V$  is the size of the vocabulary and  $i$  is the index of word  $\mathbf{x}_t$  in the vocabulary.
- We could use an MLP for sequence data by stacking the  $\mathbf{x}^{(i)}$  into a long vector:

$$\begin{bmatrix} \mathbf{x}_1^{(i)} & \mathbf{x}_2^{(i)} & \dots & \mathbf{x}_{T_i}^{(i)} \end{bmatrix}^\top \in \mathbb{R}^{VT \times 1}$$

- If the MLP has  $H$  hidden units, the weight matrix would have the dimension  $H \times VT$ .
- In practice,  $V$  can be very large (e.g.,  $V \sim 10,000$ ) in large-scale NLP tasks such as machine translation. This leads to the issue of the **curse of dimensionality**.
- Moreover, MLP treats sequence data as a *flattened vector*, losing **temporal information**.

## Language Models

- **Definition:** A probabilistic model that assigns probabilities to sequences of words. For example:

$$\mathbb{P}(\mathbf{x}_1, \dots, \mathbf{x}_T) = \mathbb{P}(\text{"The cat chases the mouse into a small hole."})$$

- Using the **chain rule**, the probability of a sequence of words can be decomposed as:

$$\begin{aligned}\mathbb{P}(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T) &= \mathbb{P}(\mathbf{x}_T \mid \mathbf{x}_1, \dots, \mathbf{x}_{T-1}) \cdot \mathbb{P}(\mathbf{x}_1, \dots, \mathbf{x}_{T-1}) \\ &= \mathbb{P}(\mathbf{x}_T \mid \mathbf{x}_1, \dots, \mathbf{x}_{T-1}) \cdot \mathbb{P}(\mathbf{x}_{T-1} \mid \mathbf{x}_1, \dots, \mathbf{x}_{T-2}) \\ &\quad \cdot \mathbb{P}(\mathbf{x}_2 \mid \mathbf{x}_1) \cdot \mathbb{P}(\mathbf{x}_1) \\ &= \mathbb{P}(\mathbf{x}_1) \cdot \prod_{i=2}^T \mathbb{P}(\mathbf{x}_i \mid \mathbf{x}_1, \dots, \mathbf{x}_{i-1})\end{aligned}$$

- For a **neural language model** (NLM), we use a neural network to model the **conditional probability** of the next word given the previous words:

$$\mathbb{P}(\mathbf{x}_{t+1} \mid \mathbf{x}_1, \dots, \mathbf{x}_t) = f_{\theta}(\mathbf{x}_1, \dots, \mathbf{x}_t),$$

where  $f_{\theta}$  is a parameterized function (e.g., a neural network) that outputs a **probability distribution** over the vocabulary for the next word.

# Recurrent Neural Networks for Language Models

To model conditional probability:

$$P(\mathbf{x}_{t+1} \mid \mathbf{x}_1, \dots, \mathbf{x}_t) = f_{\theta}(\mathbf{x}_1, \dots, \mathbf{x}_t).$$

- RNNs encode the history into a compact **hidden state**  $\mathbf{h}_t$ , i.e.,

$$(\mathbf{x}_1, \dots, \mathbf{x}_t) \mapsto \mathbf{h}_t.$$

- 

$$(\mathbf{h}_{t-1}, \mathbf{x}_t) \mapsto \mathbf{h}_t.$$

- The conditional probability is modeled through the current hidden state:

$$P(\mathbf{x}_{t+1} \mid \mathbf{x}_1, \dots, \mathbf{x}_t) = f_{\theta}(\mathbf{h}_t).$$

- Specifically, the RNN unit updates are given by:

$$\mathbf{h}_t = \tanh(\mathbf{W}_h \mathbf{h}_{t-1} + \mathbf{W}_x \mathbf{x}_t + \mathbf{b}_h),$$

$$\hat{\mathbf{y}}_t = \text{softmax}(\mathbf{W}_y \mathbf{h}_t + \mathbf{b}_y)$$

where  $\hat{\mathbf{y}}_t$  is the probability distribution over the vocabulary.

## Key Insights

RNNs capture **temporal dependencies** by updating the hidden state  $\mathbf{h}_t$  at each time step, sharing the same weights  $\mathbf{W}_h, \mathbf{W}_x, \mathbf{W}_y$  for **parameter efficiency** and **consistent** learning across sequences.



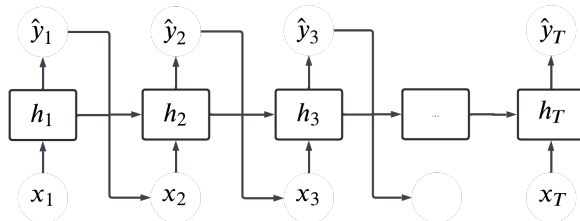
# Training RNNs

The RNN unit updates are defined as:

$$\mathbf{h}_t = \tanh(\mathbf{W}_h \mathbf{h}_{t-1} + \mathbf{W}_x \mathbf{x}_t + \mathbf{b}_h),$$

$$\hat{\mathbf{y}}_t = \text{softmax}(\mathbf{W}_y \mathbf{h}_t + \mathbf{b}_y).$$

- The training dataset is  $\mathcal{D} = \{\mathbf{y}^{(i)}\}_{i=1}^N$ , where each  $\mathbf{y}^{(i)} = \{\mathbf{y}_t^{(i)}\}_{t=1}^T$  is a sequence of words.
- Starting with  $\mathbf{x}_1 = \mathbf{y}_1$ , the model iteratively uses  $\mathbf{x}_t = \hat{\mathbf{y}}_{t-1}$  to predict  $\hat{\mathbf{y}}_t$ .

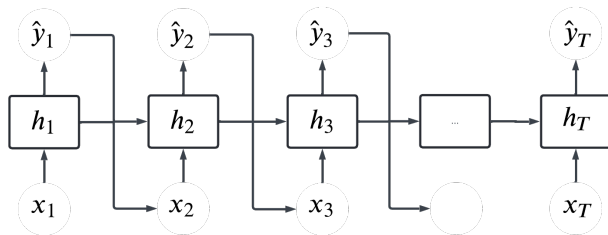


- The total cost is computed using cross-entropy loss:

$$\mathcal{L}(\theta) = -\frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T \mathbf{y}_t^{(i)} \cdot \log \hat{\mathbf{y}}_t^{(i)}$$

- **Define: Self-supervised learning** is a machine learning technique where a model generates its own labels from unlabeled data.

# Backpropagation Through Time



- **Forward propagation:**

$$\mathbf{h}_t = \phi(\mathbf{W}_h \mathbf{h}_{t-1} + \mathbf{W}_x \mathbf{x}_t), \quad \hat{\mathbf{y}}_t = \sigma(\mathbf{W}_y \mathbf{h}_t).$$

- **Backpropagation through time:**

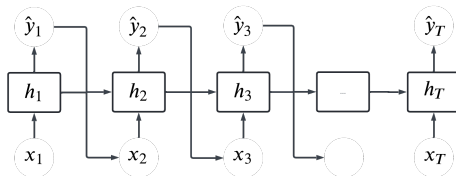
$$d\mathbf{h}_t = \mathbf{W}_h^\top [\phi'_{t+1} \odot d\mathbf{h}_{t+1}] + \mathbf{W}_y^\top [\sigma'_t \odot d\mathbf{y}_t]$$

$$d\mathbf{W}_h = \sum_t [d\mathbf{h}_t \odot \phi'_t] \mathbf{h}_{t-1}, \quad d\mathbf{W}_x = \sum_t [d\mathbf{h}_t \odot \phi'_t] \mathbf{x}_t, \quad d\mathbf{W}_y = \sum_t [d\hat{\mathbf{y}}_t \odot \sigma'_t] \mathbf{h}_t^\top$$

where  $\phi_t := \phi(\mathbf{W}_h \mathbf{h}_{t-1} + \mathbf{W}_x \mathbf{x}_t)$  and  $\sigma_t := \sigma(\mathbf{W}_y \mathbf{h}_t)$ .

# Generating a New Sequence with NLM

Once the NLM is well trained, we can sample new sequences of text by following these steps:



- 1 **Start with a seed word or token:** Choose an initial word as input.
- 2 **Feed the word to the model:** The model predicts the next word based on the history:

$$h_t = \tanh(\mathbf{W}_h h_{t-1} + \mathbf{W}_x x_t), \quad \hat{y}_t = \text{softmax}(\mathbf{W}_y h_t).$$

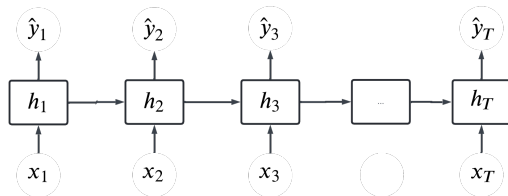
- 3 **Sample the next word:** Select the next word from the predicted probability distribution:
  - **Greedy decoding:** Pick the word with the highest probability.
  - **Stochastic sampling:** Randomly select a word based on the probability distribution.
- 4 **Iterate:** Use the **sampled** word as the input for the next time step, and repeat the process.
- 5 **Termination:** Stop when an EOS token is generated or a maximum length is reached.

**Example:** Given the current sequence “The dog chases the”, the model predicts:

$$\hat{y}_t = \{ \text{"a"} : 0.01, \text{"bird"} : 0.25, \text{"cat"} : 0.45, \dots, \text{"zoo"} : 0.01, \}$$

The model samples “cat”, feeds it back, and continues generating until the sequence ends.

# Different Types of RNNs



**One-to-Many (Sequence Generation):** A single input leads to a sequence of outputs.

- **Application:** Image captioning, music generation.

**Many-to-One (Sequence Classification):** Processes a sequence of inputs to produce a single output.

- **Application:** Sentiment analysis, speech emotion recognition, Speaker Identification/Verification.

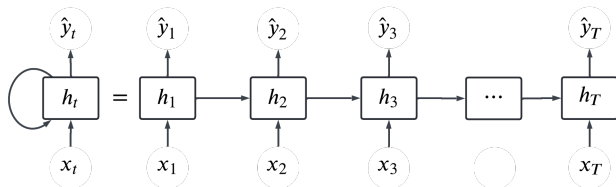
**Many-to-Many (Synchronous):** Input and output sequences have the same length.

- **Application:** Video classification, named entity recognition.

**Many-to-Many (Sequence-to-Sequence):** Input and output sequences can have different lengths.

- **Application:** Machine translation, speech recognition.

## Summary

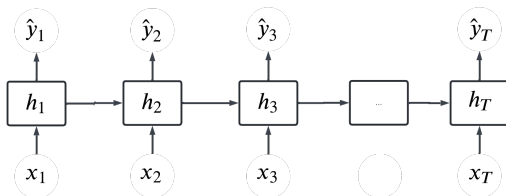


- **Challenges in Text Data:** High dimensionality and loss of temporal information.
- **Neural Language Models (NLMs):** Use neural networks to model the conditional probability of the next word given the previous ones.
- **RNNs:** Encode the history into a compact hidden state  $h_t$ , which is updated by combining the previous hidden state  $h_{t-1}$  with the current input  $x_t$ .
- **Training RNNs:**
  - Forward (simplified):  $h_t = \phi(W_h h_{t-1} + W_x x_t)$
  - Backward (simplified):  $dh_t = W_h^\top (\phi'_{t+1} \odot dh_{t+1}) + W_y^\top (\sigma'_t \odot dy_t)$
- **Generation:** Sample the next word from the predicted probability distribution produced by RNNs.
- **RNN Types:** One-to-many, many-to-one, or many-to-many structures for different tasks.

# Outline

- 1 Speech and Language Problems
- 2 Recurrent Neural Networks (RNNs)
- 3 Stabilize RNNs Learning**
- 4 Word Embedding

## Vanishing or Exploding Gradients in RNNs



- **Forward (simplified):**

$$h_t = \phi(\mathbf{W}_h h_{t-1} + \mathbf{W}_x x_t) \approx \mathbf{W}_h h_{t-1} + \mathbf{W}_x x_t \approx \mathbf{W}_h^t x_0 = \mathcal{O}(a^t)$$

- **Backward (simplified):**

$$dh_t = \mathbf{W}_h^\top (\phi'_{t+1} \odot dh_{t+1}) + \mathbf{W}_y^\top (\sigma'_t \odot dy_t) \approx \mathbf{W}_h^{\top(T-t)} dh_T = \mathcal{O}(b^{T-t})$$

- **Long-term dependencies:**

*"The dog chased a cat down the street, ran out of the house, jumped over a fence, and after running for what seemed like miles, finally caught the ..."*

# Gated Recurrent Unit (GRU)

- **Update Gate ( $z_t$ ):** Controls how much of the previous hidden state is carried forward:

$$z_t = \sigma(\mathbf{W}_z[\mathbf{h}_{t-1}, \mathbf{x}_t])$$

where  $\sigma(\cdot)$  is the sigmoid function, outputting values in the range  $(0, 1)$ .

- **Final Hidden State ( $\mathbf{h}_t$ ):** Combines the previous hidden state  $\mathbf{h}_{t-1}$  and the candidate state  $\tilde{\mathbf{h}}_t$  based on the update gate  $z_t$ :

$$\mathbf{h}_t = z_t \odot \mathbf{h}_{t-1} + (1 - z_t) \odot \tilde{\mathbf{h}}_t$$

- **Candidate Hidden State ( $\tilde{\mathbf{h}}_t$ ):** Computed by resetting parts of the previous hidden state:

$$\tilde{\mathbf{h}}_t = \tanh(\mathbf{W}_h(\mathbf{r}_t \odot \mathbf{h}_{t-1}) + \mathbf{W}_x \mathbf{x}_t)$$

- **Reset Gate ( $\mathbf{r}_t$ ):** Determines how much of the previous hidden state should be forgotten:

$$\mathbf{r}_t = \sigma(\mathbf{W}_r[\mathbf{h}_{t-1}, \mathbf{x}_t])$$

## Key Insights

- The update gate  $z_t$  helps **carry** important information from the past for long-term dependencies.
- The reset gate  $\mathbf{r}_t$  forces **discards** irrelevant past information, focusing on the most important data.



## Long Short-Term Memory (LSTM)

- **Forget Gate ( $f_t$ ):** Decides what information to discard from the previous cell state:

$$f_t = \sigma(\mathbf{W}_f[\mathbf{h}_{t-1}, \mathbf{x}_t])$$

- **Input Gate ( $i_t$ ):** Controls which new information to update in the cell state:

$$i_t = \sigma(\mathbf{W}_i[\mathbf{h}_{t-1}, \mathbf{x}_t])$$

- **Output Gate ( $o_t$ ):** Controls what part of the cell state should be output:

$$o_t = \sigma(\mathbf{W}_o[\mathbf{h}_{t-1}, \mathbf{x}_t])$$

- **Candidate Cell State ( $\tilde{c}_t$ ):** Computes the new candidate values for the cell state:

$$\tilde{c}_t = \tanh(\mathbf{W}_c[\mathbf{h}_{t-1}, \mathbf{x}_t])$$

- **Cell State ( $c_t$ ):** The cell state is updated based on the forget and input gates:

$$c_t = f_t \odot c_{t-1} + i_t \odot \tilde{c}_t$$

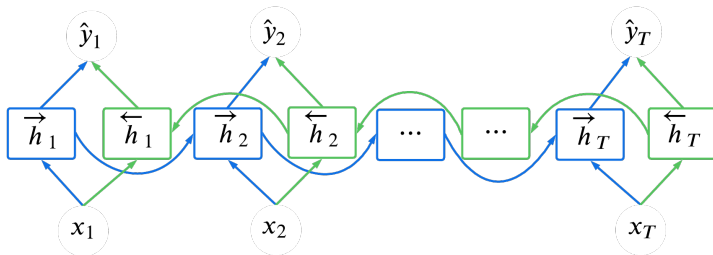
The final hidden state is:

$$\mathbf{h}_t = o_t \odot \tanh(c_t)$$

### Key Insights

- LSTMs maintain long-term dependencies via the memory cell state  $c_t$ .
- LSTMs are more flexible than GRU by using more gates to control the flow of information.

# Bidirectional Recurrent Neural Networks (BRNNs)



- **Forward:** The RNN processes the input sequence from the first time step to the last.

$$\vec{h}_t = \tanh(\mathbf{W}_f[\vec{h}_{t-1}, \mathbf{x}_t])$$

- **Backward:** Another RNN processes the sequence in reverse from the last time step to the first.

$$\overleftarrow{h}_t = \tanh(\mathbf{W}_b[\overleftarrow{h}_{t+1}, \mathbf{x}_t])$$

- **Combined Hidden State:** The final hidden state is the concatenation of the forward and backward states:

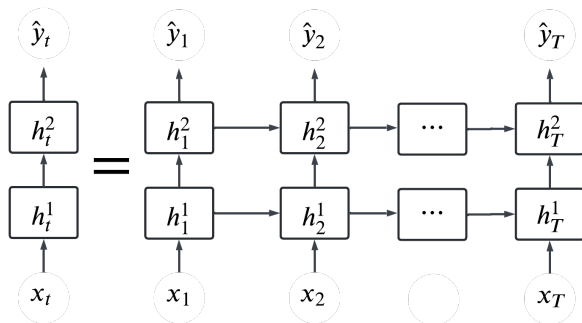
$$\mathbf{h}_t = [\vec{h}_t, \overleftarrow{h}_t]$$

## Deeper RNNs

- **Forward propagation:** Each layer  $\ell$  computes its hidden state by using the hidden state from the previous layer  $\ell - 1$ :

$$\mathbf{h}_t^{(\ell)} = \phi(\mathbf{W}_h^{(\ell)} \mathbf{h}_{t-1}^{(\ell)} + \mathbf{W}_x^{(\ell)} \mathbf{h}_t^{(\ell-1)}),$$

where  $\mathbf{h}_t^{(0)} = \mathbf{x}_t$ .



- The hidden state from the final layer is used for predictions.

# Outline

- 1 Speech and Language Problems
- 2 Recurrent Neural Networks (RNNs)
- 3 Stabilize RNNs Learning
- 4 Word Embedding**

## Featurized Representation

### Drawbacks of One-Hot Representation:

- **Orthogonality:** One-hot vectors are orthogonal, meaning they don't capture any relationships or similarities between words.
- **High Dimensionality:** One-hot vectors are sparse and grow with the size of the vocabulary, making them inefficient for large vocabularies (e.g., millions of words).

### Example: Word Correlation Matrix

Category	man	woman	king	queen	apple	orange
<b>Gender</b>	-1.00	1.00	-0.95	0.97	0.01	0.02
<b>Royalty</b>	0.01	0.02	1.00	0.9	-0.01	0.08
<b>Age</b>	0.02	0.03	0.71	0.68	0.05	0.05
<b>Food</b>	-0.01	0.02	-0.02	0.03	0.81	0.75

### Word Embedding:

- Word embeddings represent words as **dense vectors** in a **lower-dimensional space**.
- Words used in similar contexts have higher correlations, capturing their **semantic relationships**.

# Using Word Embeddings

## Embedding Matrix:

- The embedding matrix  $E \in \mathbb{R}^{d \times V}$  is a learned matrix that transforms a one-hot encoded word  $x$  into a **dense** word embedding vector  $e$ :

$$e = Ex$$

where  $d$  is the embedding dimension, and  $V$  is the size of the vocabulary.

- This transformation allows words to be represented as **dense** vectors in a **lower-dimensional** space, capturing **semantic** relationships.

## Transfer Learning:

- Word embeddings can be pre-trained on large corpora (e.g., Wikipedia or news articles), allowing models to capture rich semantic relationships.
- In small datasets, rare or specific words may be hard to learn, but pre-trained embeddings group similar words together. This helps models generalize better, even with limited data.

## Analogy:

- Word embeddings capture analogies.

$$\text{man} - \text{woman} \approx \text{king} - \text{queen}$$

- This means that models can understand not only word meanings but also deeper relationships between words.

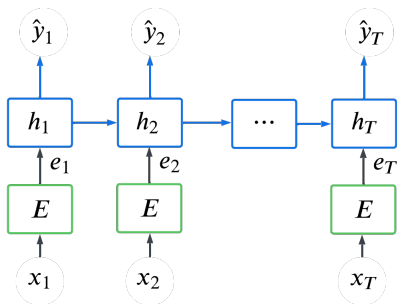
# Learning Word Embeddings with a Language Model

Let  $\{w_1, \dots, w_T\}$  be a sequence of words, e.g., “*The cat chased a mouse into a hole in the wall.*”

- In previous language models, we aim to model the probability of the next word given the history:

$$\mathbb{P}(w_t \mid w_1, \dots, w_{t-1}) = f_{\theta}(x_1, \dots, x_{t-1}),$$

where  $x_t$  is a one-hot vector representing the word at time step  $t$ .



- With word embeddings, this becomes:

$$\mathbb{P}(w_t \mid w_1, \dots, w_{t-1}) = f_{\theta}(e_1, \dots, e_{t-1}),$$

where the word embedding  $e_t$  is computed by

$$e_t = \mathbf{E}x_t,$$

$\mathbf{E} \in \mathbb{R}^{d \times V}$  is the embedding matrix, with  $d$  the embedding dimension, and  $V$  the vocabulary size.

- During training, both language model  $f_{\theta}$  and the embedding matrix  $\mathbf{E}$  are learned **simultaneously**.

## Word2vec: CBOW and Skip-Gram

**Continuous Bag of Words (CBOW):** Predicts the target word given the context:

$$\mathbb{P}(\mathbf{w}_t \mid \mathbf{w}_{t-n}, \dots, \mathbf{w}_{t+n}) = f_{\theta}(e_{t-n}, \dots, e_{t+n})$$

where the  **$n$ -gram** sequence  $\{e_{t-n}, \dots, e_{t+n}\}$  represents **context words** and  $e_t$  is the **target word**.

- Practically, a **shallow** network with  $n = \mathcal{O}(1)$  is sufficient to learn embedding matrix  $\mathbf{E}$ .

$$e_{t \pm i} = \mathbf{E}x_{t \pm i}, \quad \mathbf{h} = \mathbf{W}[e_{t-n}, \dots, e_{t+n}], \quad \hat{\mathbf{y}} = \text{softmax}(\mathbf{h}), \quad \forall i \in \{1, 2, \dots, n\},$$

where the concatenation  $[e_{t-n}, \dots, e_{t+n}]$  can be replaced with summation or averaging to reduce computational complexity.

**Skip-Gram:** Predicts the context words given a target word.

$$\mathbb{P}(\mathbf{w}_{t-n}, \dots, \mathbf{w}_{t+n} \mid \mathbf{w}_t) = f_{\theta}(e_t)$$

- For each context word  $\mathbf{w}_{t \pm i}$ , we predict its probability given the embedding of target word  $e_t$ :

$$e_t = \mathbf{E}x_t, \quad \mathbf{h}_i = \mathbf{W}e_t, \quad \hat{\mathbf{y}}_{t \pm i} = \text{softmax}(\mathbf{h}_i), \quad \forall i \in \{1, 2, \dots, n\}.$$



# Negative Sampling

- **Expensive softmax:** The full softmax is computationally expensive because it requires summing over the entire vocabulary:

$$\text{softmax}(\mathbf{h}) = \frac{e^{\mathbf{h}_j}}{\sum_{i=1}^V e^{\mathbf{h}_i}}$$

where  $V \sim 1$  million.

- Reformulates the context-target prediction as a **binary** classification:  $(\mathbf{w}_t, \mathbf{w}_c, y)$ , where  $y$  is label.
- **Binary classifier for context-target pair:**  $\mathbb{P}(y = 1 \mid \mathbf{w}_c, \mathbf{w}_t) = \sigma(\mathbf{e}_c^\top \mathbf{e}_t)$ , where  $\sigma(\cdot)$  is sigmoid
- **Maximize log-likelihood:** It is equivalent to minimizing the following objective:

$$\mathcal{L}(\mathbf{E}) = - \sum_{(\mathbf{e}_c, \mathbf{e}_t)} \log \sigma(\mathbf{e}_t^\top \mathbf{e}_c) - \sum_{(\mathbf{e}_c, \tilde{\mathbf{e}}_t)} \log \sigma(-\mathbf{e}_t^\top \tilde{\mathbf{e}}_t)$$

where  $\tilde{\mathbf{e}}_t$  are **negative** samples from words outside the context window.

# Word Analogy Task

- **Objective:** Assess the quality of word embeddings by testing how well they capture **semantic** and **syntactic** relationships between words.
- **Goal:** Given an analogy of the form: “*a is to b as c is to ???*”, find the word *d* that completes the analogy correctly.
- **Examples:**
  - Semantic analogy: “*Paris is to France as Washington, D.C. is to the United States.*”
  - Syntactic analogy: “*Run is to Running as Swim is to Swimming.*”
- **Vector Arithmetic:**  $e_d \approx e_b - e_a + e_c$ , i.e., *linear space*
- The word *d* is chosen as the vector  $e_d$  closest to  $e_b - e_a + e_c$  based on *cosine similarity*:

$$d = \arg \max_{x \in V} S_C(e_x, e_b - e_a + e_c)$$

where  $V$  is the vocabulary, and

$$S_C(\mathbf{u}, \mathbf{v}) = \cos(\theta) = \frac{\mathbf{u}^\top \mathbf{v}}{\|\mathbf{u}\| \cdot \|\mathbf{v}\|}$$

where  $\theta$  is the angle between vectors  $\mathbf{u}$  and  $\mathbf{v}$ .

- **Evaluation:** The predicted word *d* is evaluated by comparing it to the correct answer from the analogy dataset.

# GloVe: Global Vectors for Word Representation

- **Objective:** Create a word embedding model that captures both local context and **global** statistical information from a text corpus.

- **Co-occurrence Matrix:**

- $X_{ij}$ : Number of times word  $j$  appears in the context of word  $i$ .
- $X_i = \sum_j X_{ij}$ : Total occurrences of any word in the context of word  $i$ .
- $\mathbb{P}(w_j | w_i) = X_{ij}/X_i$ : Probability of word  $j$  occurs in the context of word  $i$ .

- **Word Comparison in Context:**

- Compare words  $e_i$  and  $e_j$  in the context  $\tilde{e}_k$  using a **probability ratio**:

$$\exp \left\{ (e_i - e_j)^\top \tilde{e}_k \right\} = \frac{\mathbb{P}(w_i | w_k)}{\mathbb{P}(w_j | w_k)} = \frac{X_{ki}}{X_{kj}}$$

- Taking the log yields:

$$e_i^\top \tilde{e}_k - e_j^\top \tilde{e}_k = \log X_{ki} - \log X_{kj} \quad \Rightarrow \quad e_i^\top \tilde{e}_k \sim \log X_{ki}$$

- **Cost Function:**

- The GloVe model learns word vectors  $e_i$  and context vectors  $\tilde{e}_k$  by minimizing:

$$J = \sum_{i,k} f(X_{ik}) \left( e_i^\top \tilde{e}_k + b_i + \tilde{b}_k - \log(X_{ik}) \right)^2$$

- $f(X_{ik})$  is a weighting function for co-occurrences, while  $b_i$  and  $\tilde{b}_k$  are bias to maintain symmetry.

## Bias in Word Embeddings

- **Problem:** Word embeddings trained on large datasets often encode societal biases, like gender stereotypes:

*“Man is to Computer Programmer as Woman is to Homemaker?”*

- **Identifying Bias Direction:** Compute the average difference between **definitional pairs**:

$$\begin{cases} \vec{e}_{\text{man}} - \vec{e}_{\text{woman}} \\ \vec{e}_{\text{he}} - \vec{e}_{\text{she}} \\ \dots \end{cases} \Rightarrow \vec{e}_{\text{bias}} = \frac{1}{N} \sum_{i=1}^N (\vec{e}_{x_i} - \vec{e}_{y_i})$$

This average can be replaced by advanced techniques like PCA.

- **Neutralization:** Project non-definitional words onto the space orthogonal to the bias direction to remove bias:

$$\vec{e} \leftarrow \vec{e} - \frac{\vec{e}^T \vec{e}_{\text{bias}}}{\|\vec{e}_{\text{bias}}\|^2} \cdot \vec{e}_{\text{bias}}$$

- **Equalizing Pairs:** Adjust equalize pairs (like “brother” and “sister”) to have equal and opposite projections along the gender direction, making them **equidistant** from the gender-neutral axis.
- **Identifying Gendered Words:** Train a classifier to distinguish between gender-specific and neutral words using a set of definitional pairs.