# CS 333:

# Natural Language Processing

Fall 2025

Prof. Carolyn Anderson Wellesley College

# Recent Work by Wellesley Alums

# ReasoningWeekly: A General Knowledge and Verbal Reasoning Challenge for Large Language Models

Zixuan Wu, Francesca Lucchetti, Aleksander Boruch-Gruszecki, Jingmiao Zhao, Carolynane Anderson, Joydeep Biswas, Federico Cassano, Arjun Guha

- July, Seni r Area Chairs, Area Chairs, Reviewers, Authors, Commitment Readers
  Revisions
- © CC BY 4.0

#### Abstract:

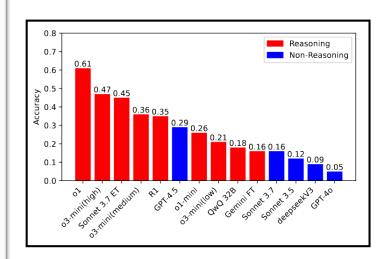
Existing bench marks for frontier models often test specialized, "PhD-level" knowledge that is difficult for non-experts to grasp. In contrast, we present a benchmark with 613 problems based on the NPR Sunday Puzzle Challenge that requires only general knowledge. Our benchmark is challenging for both humans and models; however correct solutions are easy to verify, and models' mistakes are easy to spot. As LLMs are more widely deployed in society, we believe it is useful to develop benchmarks for frontier models that humans can understand without the need for deep domain expertise.

Class of 2024



Class of 2023 (took CS 333 in Fall 2023!)





#### **Example item:**

Challenge: Think of a common greeting in another a country that is not the United States. You can rearrange its letters to get the capital of a country that neighbors the country where this greeting is commonly spoken. What greeting is it?

Ground Truth Answer: Ni hao -> Hanoi

Accepted to the International Joint Conference on Natural Language Processing & Asia-Pacific Chapter of the Association for Computational Linguistics, 2025 (AACL)

# Major NLP Conferences

- main votating conference ACL Asio AACL Nations of America NAACL - Europe ACL SACL - Empirical Methods in NLP EMNLP Conference on Longrage Modeling Colm -ALL Anmology Neurlps ML: ICLR

# Reminders

- \* Midterm 2 is on Friday!
  - The exam is open-note but closed-device
  - List of topics posted on the course website.
- Tuesday -> Monday cycle for final two assignments:
  - HW 7 will be released Friday but not due until Monday, 11/17
  - HW 8 will be released on Tuesday, 11/18 and due on Monday, 11/24
- My next help hours: Thursday 4-5

#### New Research in Religious Studies

# Hinduism Online: Mechanized Murti and Automated Adoration

Presented by Dr. Holly Walters

Department of Anthropology



SCAN QR CODE TO RSVP



Tuesday, November 4th at 5pm FND 120



Recorded?

exp: 11/5
?: mk110@wellesley.edu



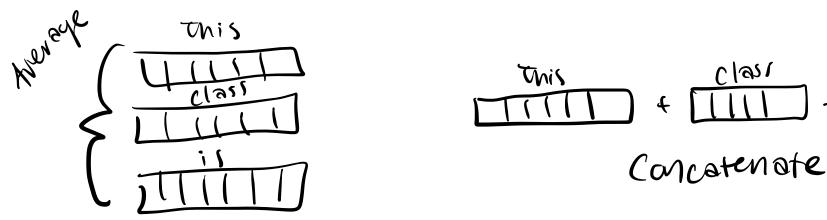
Lindsey Cameron Colloquium

November 14th, 3:45-5pm in H105

**Title:** Scalable Subjugation: The Myth of Geographic Scalability in the Gig Economy and How Workers Reconstitute Platforms

# HW 5 Curiosity Points

- Model changes and tweaks:
  - Implemented weighting to address class imbalance
  - Exploration of lexical diversity-based features
  - Extension to a neural network model
  - Visualization of feature correlations
  - Data analysis functions to aid in feature engineering
- Task extensions:
  - Experiments on different artists
  - Naive Bayes comparison experiment
  - Binary classification experiment
  - Swift album classifier
- Research literature exploration



# Recap: Recurrent Neural Networks

# **A RNN Language Model**

#### output distribution

$$\hat{y} = \text{softmax}(W_2 h^{(t)})$$

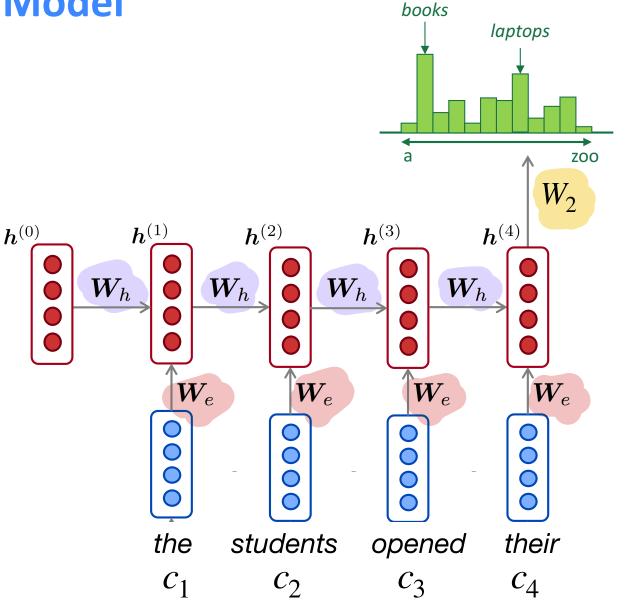
#### hidden states

$$h^{(t)} = f(W_h h^{(t-1)} + W_e c_t)$$

h<sup>(0)</sup> is initial hidden state!

#### word embeddings

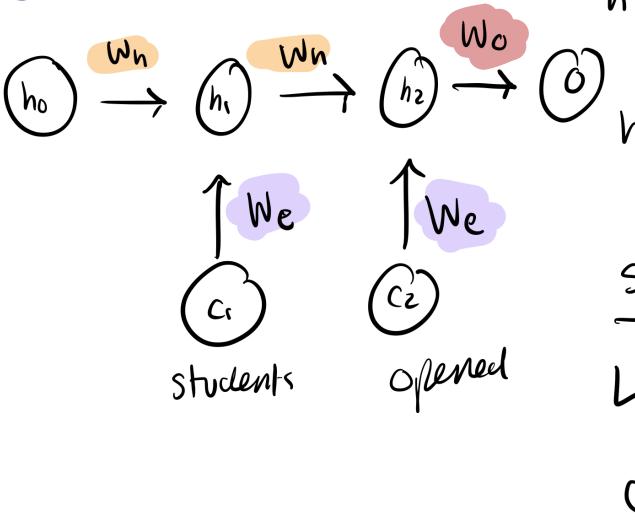
$$c_1, c_2, c_3, c_4$$



 $\hat{\boldsymbol{y}}^{(4)} = P(\boldsymbol{x}^{(5)}|\text{the students opened their})$ 

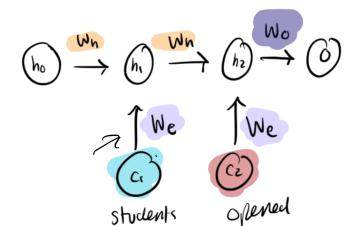
# Training a Recurrent Neural Network

# Training an RNN



# Key Question: what are the parameters?

### **Gradients**



$$\frac{\partial L}{\partial w_0} = \frac{\partial L}{\partial w_0} \cdot \frac{\partial w}{\partial w_0} = -(y-0) \cdot h_2$$

$$-(y-0) \cdot h_2$$

## Gradient of c2

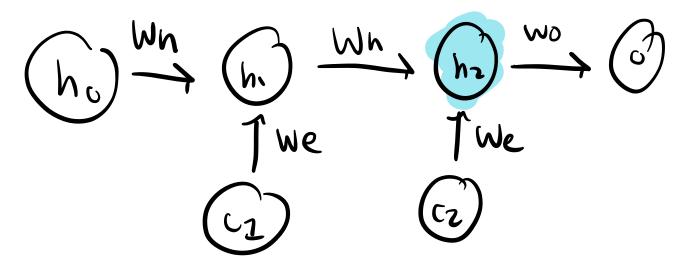
# Gradient of we

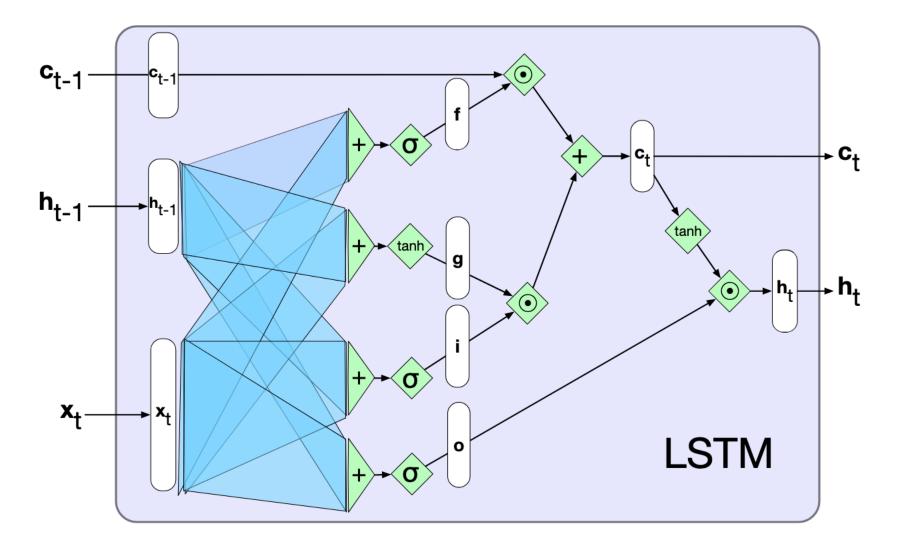
$$h_z = fanh(a+b)$$
 $a = W_e C_z$ 
 $b = W_h \cdot h_1$ 
 $h_1 = fanh(W_e C_1 + W_h \cdot h_0^2)$ 
 $d_z = W_e \cdot C_1$ 

$$\frac{\partial L}{\partial we} = \frac{1}{100} \frac{\partial L}{\partial w} \frac{\partial$$

# Vanishing Gradients Problem

- It's very easy for the gradients in an RNN to zero out since the dependencies between them are so complex.
- There are very few "pipes" for the information to flow through.





#### A single LSTM unit displayed as a computation graph.

**Inputs:** the current input, x, the previous hidden state,  $h_{t-1}$ , and the previous context,  $c_{t-1}$ . **Outputs**: a new hidden state,  $h_t$  and an updated context,  $c_t$ .

#### **Components:**

- add gate/input gate (i): selects the information to add to the current context
- forget gate (f): delete information from the context that is no longer needed
- output gate (o): decides what information is required for the current hidden state

# RNN Advantages and Limitations

#### $\hat{\boldsymbol{y}}^{(4)} = P(\boldsymbol{x}^{(5)}|\text{the students opened their})$

## why is this good?

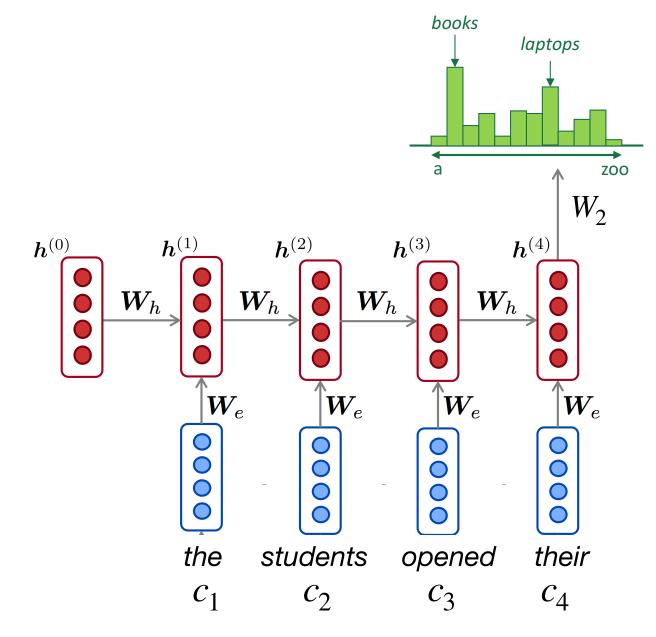
#### RNN Advantages:

- Can process any length input
- Model size doesn't increase for longer input
- Computation for step t can (in theory) use information from many steps back
- Weights are shared across timesteps > representations are shared

#### RNN **Disadvantages**:

- Recurrent computation is slow
- In practice, difficult to access information from

\_\_many steps back



## RNNs suffer from a **bottleneck** problem

 $oldsymbol{h}^{(1)}$ 

 $W_e$ 

the

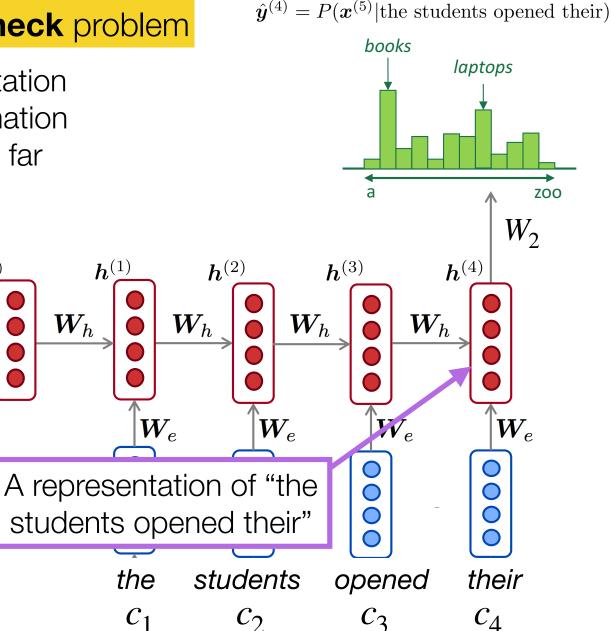
 $c_1$ 

 $oldsymbol{W}_h$  ,

 $h^{(0)}$ 

The current hidden representation must encode all of the information about the text observed so far

This becomes difficult especially with longer sequences

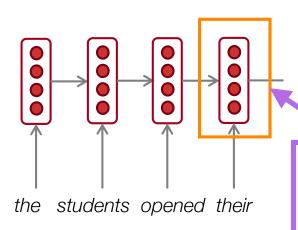


 $C_4$ 

# "you can't cram the meaning of a whole %&@#&ing sentence into a single \$\*(&@ing vector!"

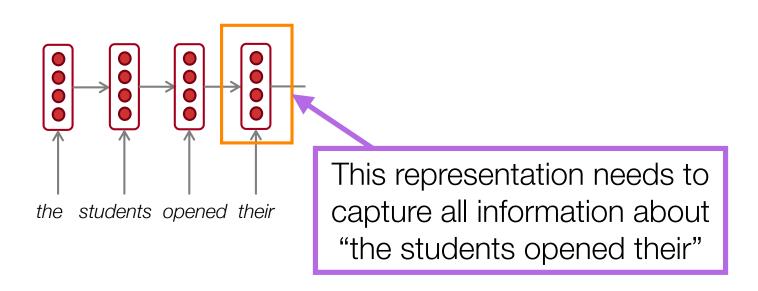
Ray Mooney (NLP professor at UT Austin)

# idea: what if we use multiple vectors?



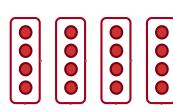
This representation needs to capture all information about "the students opened their"

# idea: what if we use multiple vectors?



## Instead of this, let's try:

the students opened their =

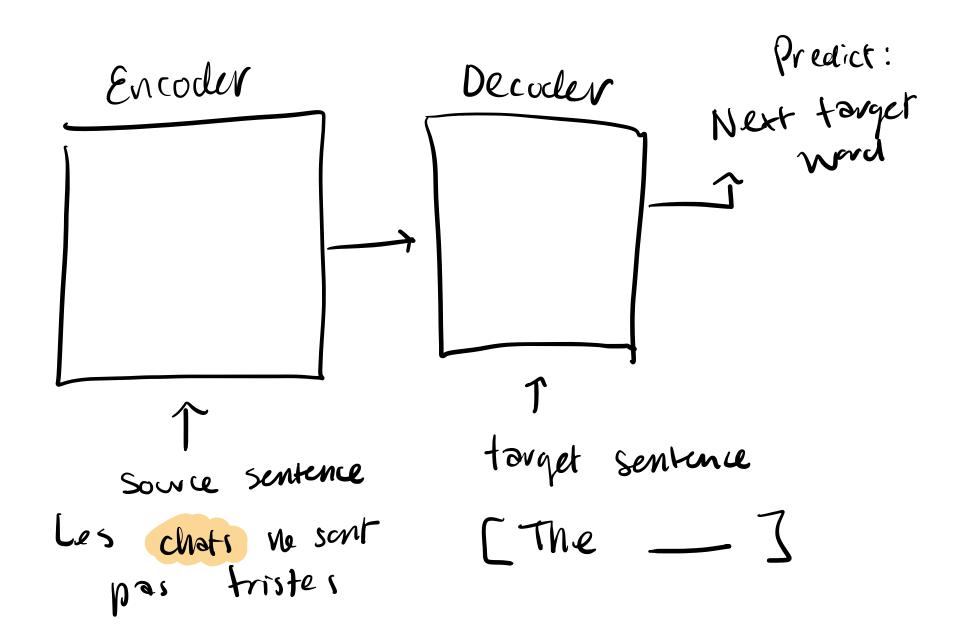


(all 4 hidden states!)

# The solution: attention

- Attention mechanisms (Bahdanau et al., 2015) allow language models to focus on a particular part of the observed context at each time step
  - Originally developed for machine translation, and intuitively similar to word alignments between different languages

## **Encoder-Decoder Models of Machine Translation**



# Attention

# How does it work?

 in general, we have a single query vector and multiple key vectors. We want to score each query-key pair

in a neural language model, what are the queries and keys?

# What Is Attention?

~	tosk-specific searching for in from the casure the sim	mputant through	<b>7.</b>
-3.4 1, = K1.9	2.4 12= K2'9	$-0.8$ $v_3 = k_3 \cdot 9$	-1.2 14 = £4.7
The	Stydats	gened	
Keys: representations of past contex			

## What Is Attention?

Guery: as a tosk-specific Vector that is "searching" for imputant things from the post context.

Softmax

Step 2: Scale scores to between 0-1

0.01 0.7 0.24 0.05
$$-3.4 \quad 2.4 \quad -0.8 \quad -1.2$$

$$1_1 = k_1 \cdot 9 \quad 1_2 = k_3 \cdot 9 \quad 1_4 = k_4 \cdot 1$$
The Studies gened than
$$Keys: Vernesent attacks of past contex$$

# What Is Attention?

as a tosk-specific vector that is Guery: "searching" for important things from the post context. Compte à weighted average of Step 3: Other: a vector of the same divensions 0.7 2.4 0.01 -0.8 -3.4 13 = k3.9 14 = k4.9 12 = k2'9 r, = k, 9 Studenti gened then The Keys: representations of



They don't tell you this in the paper (well they do but you have to read it like 15 times)



Multiplying
a lot of vectors
a lot of times
with scaled softmax

Attention