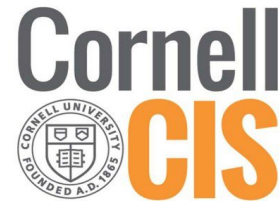




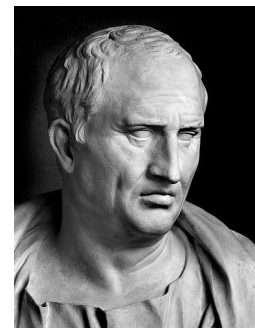
- Use Artificial Intelligence and Machine Learning methods to model human language
- To build a system that generates quotations which are indistinguishable from human quotations

- **Solutions**

- N-gram probabilistic Language Models
- Recurrent Neural Net based Language Model
- Long Short Term Memory and Gated Recurrent Unit based Language Models



- First Iteration: 186 unique quotations pulled exclusively from BrainyQuote
- Second Iteration: 4,354 unique sentences from Republic of Cicero, Treatises on Friendship and Old Age, Letters of Marcus Tullius Cicero, and a variety of transcribed speeches to augment initial dataset



- First Iteration: 577 quotations pulled exclusively from BrainyQuote
- Second Iteration: 2,379 tweets from Donald Trump's official twitter account to replace initial dataset





- $P(w_i | w_{i-1}) = \frac{\text{count}(w_i, w_{i-1})}{\text{count}(w_{i-1})}$
- Compute the word frequency of the current word given the immediately preceding word

- **Trigram Language Model**

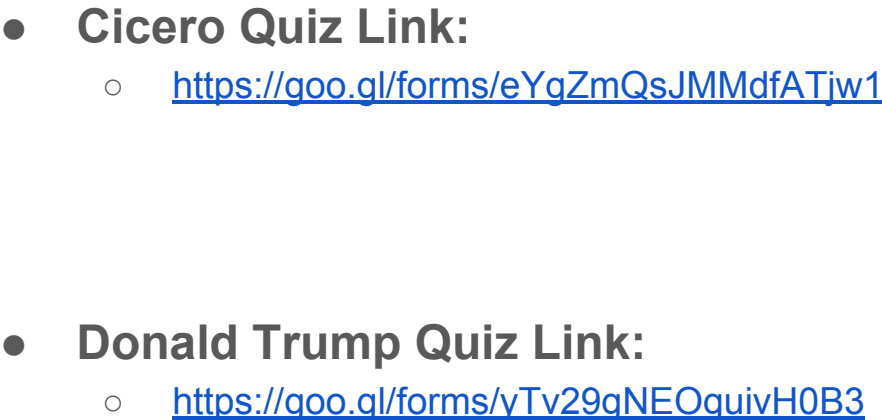
- $P(w_i | w_{i-1}, w_{i-2}) = \text{count}(w_i, w_{i-1}, w_{i-2}) / \text{count}(w_{i-1}, w_{i-2})$
- Compute the word frequency of the current word given the past two words

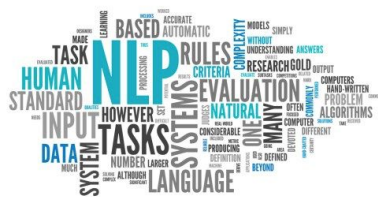
- **Handling Unseen N-grams**

- Smoothing Techniques
 - Add one
 - K-smoothing
 - Kneser Ney

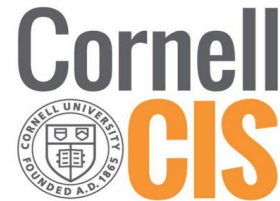
- **Handling Unknown Words**

- Introduce a 'UNK' token into vocabulary
- Map all rare words (based on word frequency) to this token
- Put uniform random distribution over all words mapped to 'UNK' token





Recurrent Neural Net Language Model



- **Intuition**

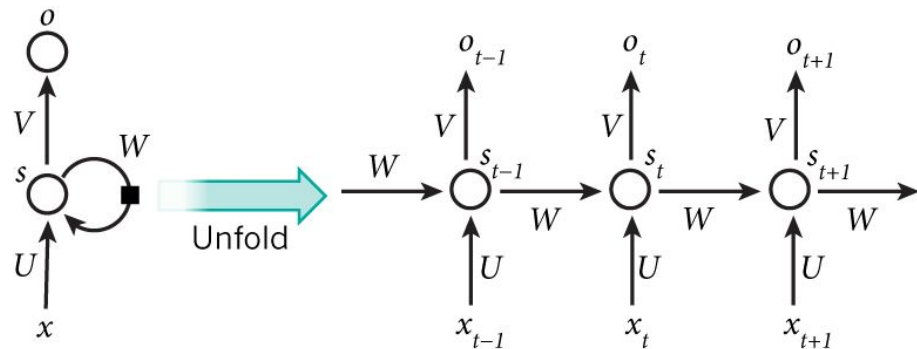
- Look at the entire past sentence when predicting the next word
 - *I grew up in France, so I am fluent in _____*
 - Earlier location (i.e. France) should influence a later language (i.e. French)

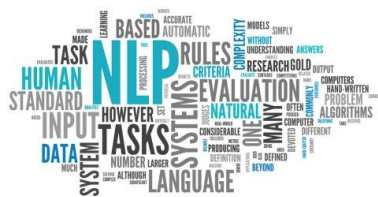
- **Numpy Implementation**

- Following a tutorial to help with the math, built a Numpy RNN from the ground up
 - Poor results on small dataset
 - Incredibly slow on larger dataset

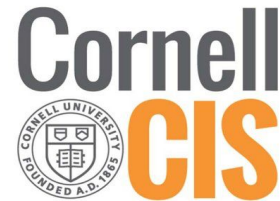
- **Theano Implementation**

- Uses GPU for matrix operations
- Poor results on large dataset



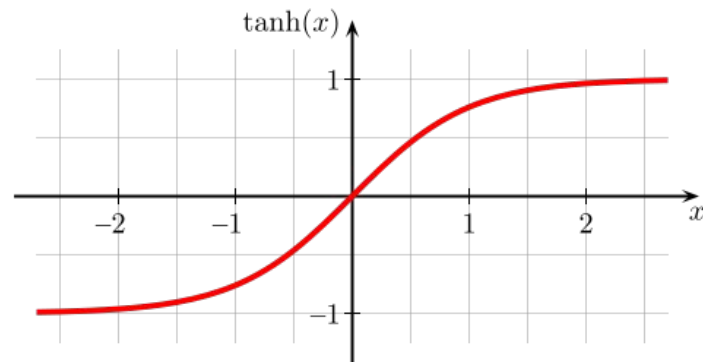


Recurrent Neural Net Language Model



- **Vanishing Gradient**

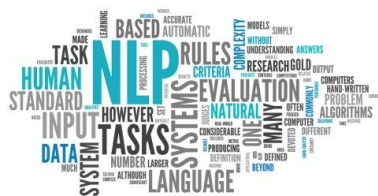
- Gradient of tanh activation function at tails is 0
- Saturated neuron in unrolled chain drives gradient down for past time steps
- Longer the unrolled chain, the more effect
- Particularly relevant for language models: length of the chain (i.e. depth) dependent on sentence length



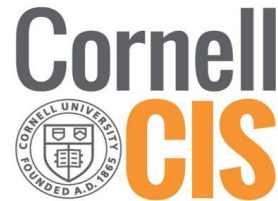
- **Effect**

- RNN unable to learn long or even mid-range inter-sentence word dependencies





Recurrent Neural Net Language Model



- **Cicero Generated Sentences**

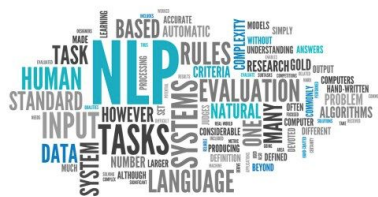
- Most RNN based Language Model sentences are poorly formed. Best results occur with short sentences.
 - “Messengers from men is appreciated seeing.”
 - “Describe the roman difficult fire that receiver demanding.”
 - “Lifetime my republic he beware the surprised method.”

- **Trump Generated Tweets**

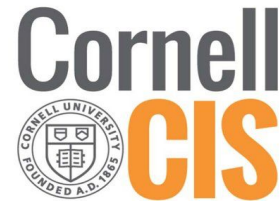
- Inherently chaotic structure of tweet => leniency for generated sentence
 - twitter good!
 - laughingstock president food.
 - shots very fired continue liar and party disaster @ turbines change.

- **Analysis**

- Many poorly formed sentences
- Unable to capture natural sentence structure (i.e. subject - verb - object)



Long Short Term Memory Language Model

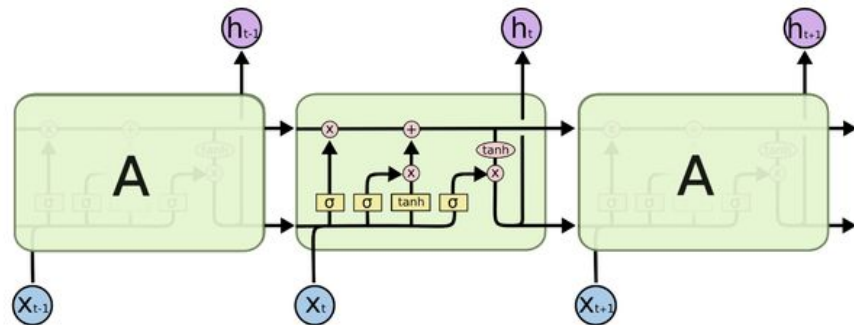


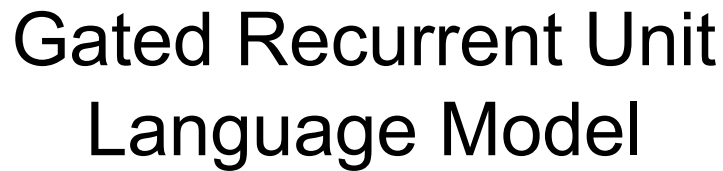
- **LSTM Architecture**

- Developed in 1997 to solve the vanishing gradient issue
- Internal Memory in addition to hidden state at each timestep

- **Gating Mechanisms**

- Input gate: creates a candidate state with input
- Forget gate: what information from past time step should propagate to current time step
- Output gate: determines final hidden state



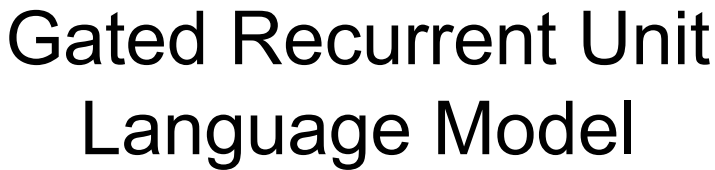


- **Gating Mechanisms**

- **Embedding Layer**

-

$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$



- For themselves I will say against to the are beginning of a will I either discharge this in the concerning up in virtue.
- What honour roman be senate to which titles ?
- I every do not consul.

- **Trump Generated Tweets**

- I want \$ in calls .
- president obama recognition - sad!
- sadly they are not smart?

- **Analysis**

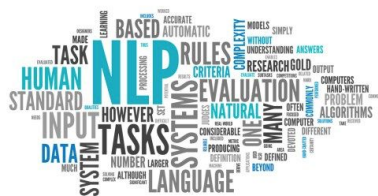
- Many poorly formed sentences
- The longer the sentence, the more likely something will go wrong
- Learns the correct inter-sentence word pattern: subject-verb, verb-object, etc.



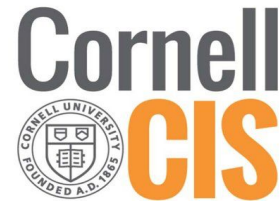
- # Trump Corpus

Cicero Corpus

Language Model	Perplexity
Bigram	358.65
Trigram	139.64
Recurrent Neural Net (Theano)	773.08
GRU Neural Net (Theano)	116.85



Extrinsic Evaluation

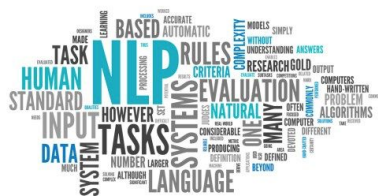


- **Cicero Quiz Results:**

- First iteration BrainyQuote dataset
- 21 respondents
- Bigram model: 66.3% average correctness score
- Trigram model: 53.9% average correctness score
- Average total scores higher than expected score ($p < 0.001$)
- Quiz takers had more trouble with trigram model vs. bigram model ($p = 0.0001186$)

- **Trump Quiz Results:**

- First Iteration BrainyQuote dataset
- 24 respondents
- Bigram model: 62.8% average correctness score
- Trigram model: 60.7% average correctness score
- Average total scores higher than expected score ($p < 0.001$)



Final Remarks



- **Overfitting Can Be Good!**

- N-gram models trained on small corpus of ‘quality’ data produced amazingly good results
- Amalgamate partial sentences from training data to produce quotation
- Overfitting to training data is good. We want to produce quotations which are likely to have been in the training data

- **GRU does great on tweets**

- Added complexity of GRU based language model
- Inherent chaotic nature and malleable inter-sentence word ordering
- Suprasses n-gram language model in terms of tweets produced
- Doesn't piece sentences together (as in the n-gram models) but learns inherent tweet structure (i.e. following a verb, a noun will have a high probability of being selected).