

# Sequential Language Modeling

Northwestern EECS 395/495 Probabilistic Graphical Models

Fall 2014

# Language Modeling

---

- ▶ **Modeling sequences of words (today)**
  - ▶ N-gram Models
  - ▶ HMMs (briefly)
  - ▶ Neural Network Language Models
  
- ▶ **Modeling Documents (next time)**
  - ▶ “Bag of words”
  - ▶ Latent Semantic Analysis, Latent Dirichlet Allocation



# Modeling Sequences of Words

---

- ▶ Statistical language models assign **probabilities** to **sequences of words**

$$P(\text{"the dog barked"}) = 4.203 * 10^{-9}$$

- ▶ Applications
  - ▶ Speech Recognition
  - ▶ Machine Translation
  - ▶ Spelling Correction
  - ▶ Information Extraction



# Outline

---

- ▶ **N-gram Models**
  - ▶ Sparsity: Smoothing, Backoff
  - ▶ Perplexity Measure
- ▶ **Exploiting Word Similarity**
  - ▶ Brown Clustering
  - ▶ HMMs (briefly)
  - ▶ Neural Network Language Models
- ▶ **Future Directions**



# So you want to Learn an n-gram model

---

- ▶ Don't implement one yourself

- ▶ SRILM

- <http://www.speech.sri.com/projects/srilm/>



# Outline

---

- ▶ **N-gram Models**
  - ▶ Sparsity: Smoothing, Backoff
  - ▶ **Perplexity Measure**
- ▶ **Exploiting Word Similarity**
  - ▶ Brown Clustering
  - ▶ HMMs (briefly)
  - ▶ Neural Network Language Models
- ▶ **Future Directions**



# How to measure LM performance?

---

- ▶ Test model  $M$  on held-out text  $\bar{w}$  of length  $N$

$$\text{Perplexity}(M, \bar{w}) = 2^{\frac{-\log_2 P_M(\bar{w})}{N}}$$



Models trained on 38 million words from the Wall Street Journal (WSJ) using a 19,979 word vocabulary.

Evaluate on a disjoint set of 1.5 million WSJ words

Perplexity: Trigram < Bigram < Unigram (in this case)

Which N-gram is the best in general?

	Unigram	Bigram	Trigram
Perplexity	962	170	109

<http://www.cs.stonybrook.edu/~ychoi/cse628/lecture/02-ngram.pdf>





### **Unigram**

Months the my and issue of year foreign new exchange's september were recession ex-  
change new endorsed a acquire to six executives

### **Bigram**

Last December through the way to preserve the Hudson corporation N. B. E. C. Taylor  
would seem to complete the major central planners one point five percent of U. S. E. has  
already old M. X. corporation of living on information such as more frequently fishing to  
keep her

### **Trigram**

They also point to ninety nine point six billion dollars from two hundred four oh six three  
percent of the rates of interest stores as Mexico and Brazil on market conditions



# Outline

---

- ▶ N-gram Models
  - ▶ Sparsity: Smoothing, Backoff
  - ▶ Perplexity Measure
- ▶ **Exploiting Word Similarity**
  - ▶ **Brown Clustering**
  - ▶ HMMs (briefly)
  - ▶ Neural Network Language Models
- ▶ **Future Directions**



# motivation

---

Corpus includes:

“Deadline on Monday”

“Deadline on Thursday”

... suggests  $P(\text{“Deadline on \_\_\_\_”})$  ?



# Brown Clustering

---

Say  $C : \mathcal{V} \rightarrow \{1, 2, \dots, k\}$  is a *partition* of the vocabulary into  $k$  classes

The model:

$$p(w_1, w_2, \dots, w_T) = \prod_{i=1}^n p(w_i | C(w_i)) p(C(w_i) | C(w_{i-1}))$$

Michael Collins

<http://www.cs.columbia.edu/~cs4705/lectures/brown.pdf>

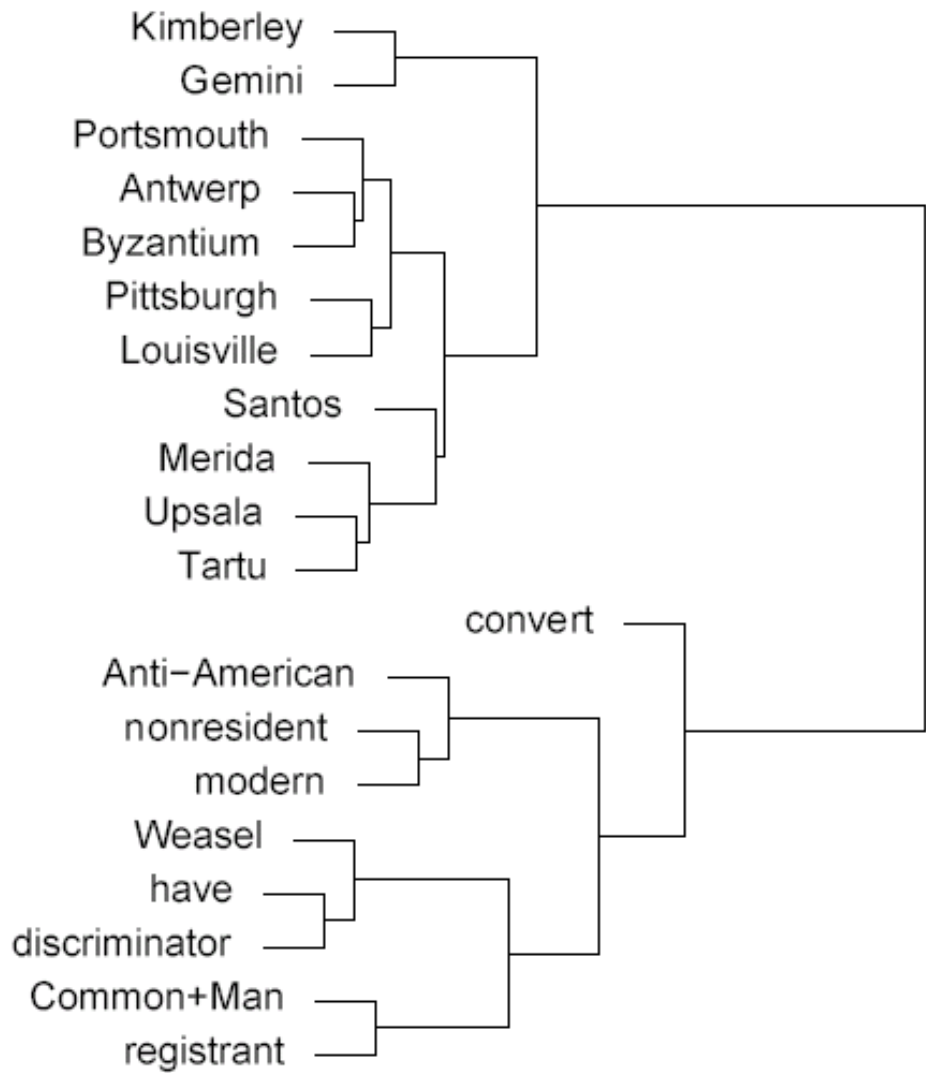


# Learning Clusters?

---

- ▶ **Greedy Agglomerative Clustering**
  - ▶ Start with each word in own cluster, iteratively combine the two clusters that results in the smallest decrease in likelihood
  - ▶ Slow (cubic in vocab size), faster approaches exist





lawyer	1000001101000
newspaperman	100000110100100
stewardess	100000110100101
toxicologist	10000011010011
slang	1000001101010
babysitter	100000110101100
conspirator	1000001101011010
womanizer	1000001101011011
mailman	10000011010111
salesman	100000110110000
bookkeeper	1000001101100010
troubleshooter	10000011011000110
bouncer	10000011011000111
technician	1000001101100100
janitor	1000001101100101
saleswoman	1000001101100110

Miller, Scott, Jethran Guinness, and Alex Zamanian. "Name Tagging with Word Clusters and Discriminative Training." *HLT-NAACL*. Vol. 4. 2004.



Nike	1011011100100101011100
Maytag	10110111001001010111010
Generali	10110111001001010111011
Gap	1011011100100101011110
Harley-Davidson	10110111001001010111110
Enfield	101101110010010101111110
genus	101101110010010101111111
Microsoft	10110111001001011000
Ventritex	101101110010010110010
Tractebel	1011011100100101100110
Synopsys	1011011100100101100111
WordPerfect	1011011100100101101000





John	101110010000000000
Consuelo	101110010000000001
Jeffrey	101110010000000010
Kenneth	10111001000000001100
Phillip	101110010000000011010
WILLIAM	101110010000000011011
Timothy	10111001000000001110
Terrence	101110010000000011110
Jerald	101110010000000011111
Harold	1011100100000000100
Frederic	1011100100000000101
Wendell	101110010000000011

# HMMs

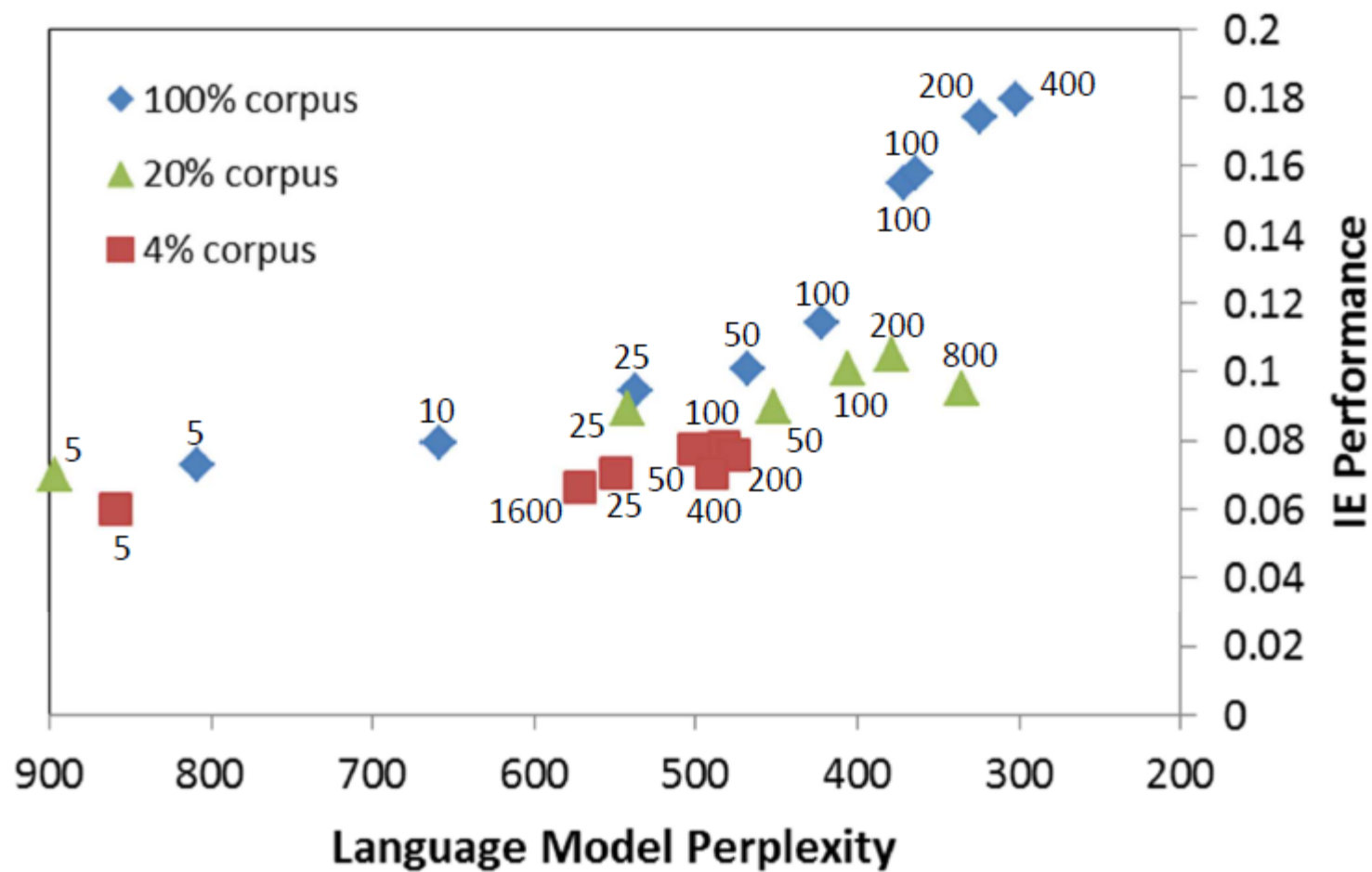
---

- ▶ “Soft” version of Brown Clusters
  - ▶ as GMMs are to K-means

model	AUC
I-HMM-TYPE-R	<b>0.18</b>
HMM-TYPE-R	0.17
BROWN-TYPE-R-3200	0.16
BROWN-TYPE-R-1000	<b>0.18</b>
BROWN-TYPE-R-320	0.15
BROWN-TYPE-R-100	0.13
LATTICE-TYPE-R	0.11
N-GRAM-R baseline	0.10
Random baseline	0.10

---



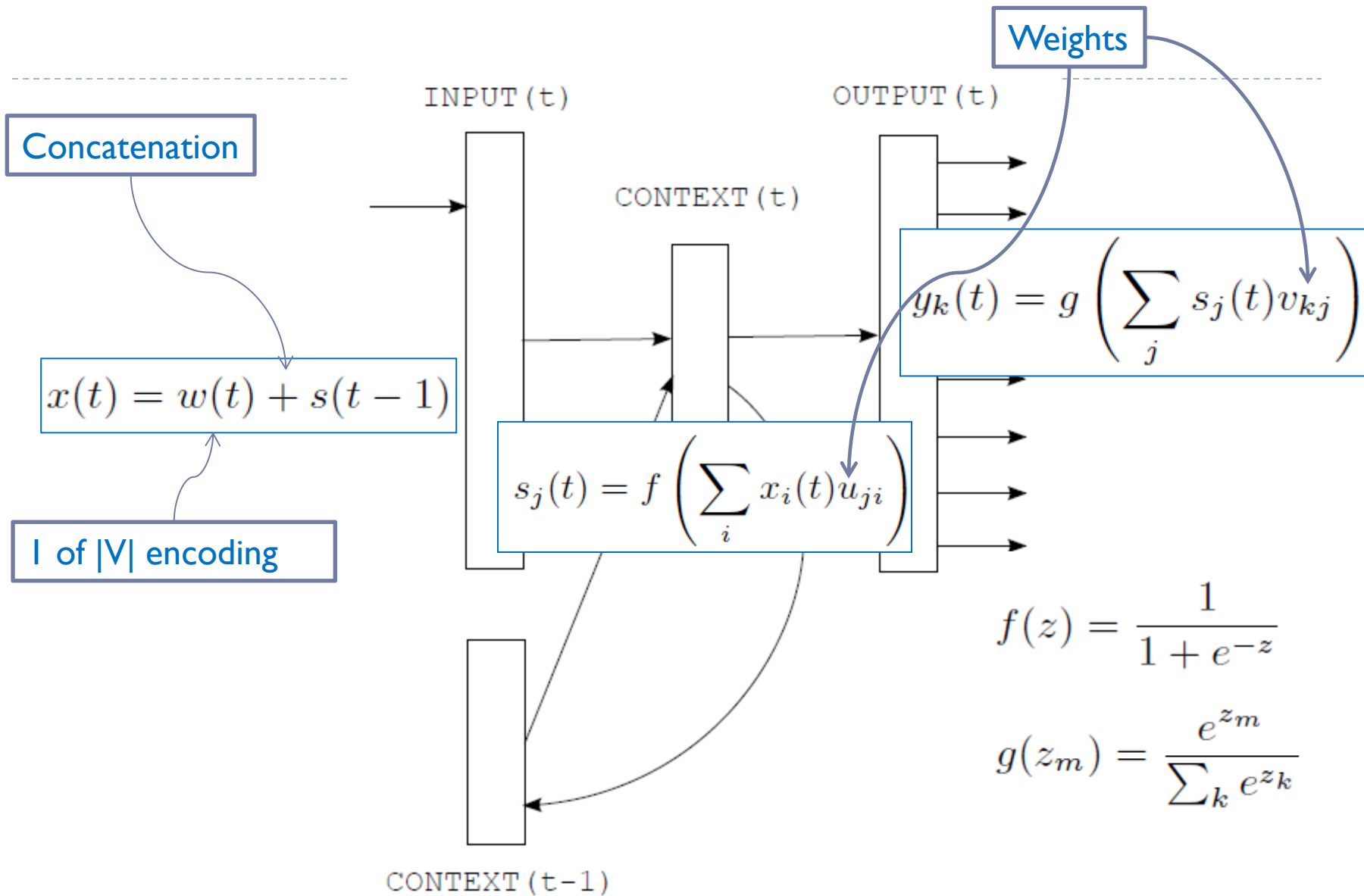


# Outline

---

- ▶ N-gram Models
  - ▶ Sparsity: Smoothing, Backoff
  - ▶ Perplexity Measure
- ▶ Exploiting Word Similarity
  - ▶ Brown Clustering
  - ▶ HMMs (briefly)
  - ▶ **Neural Network Language Models**
- ▶ **Future Directions**





# Results

---

$$0.75 * \text{RNN} + 0.25 * \text{KN}$$


Model	PPL		WER	
	RNN	RNN+KN	RNN	RNN+KN
KN5 - baseline	-	221	-	13.5
RNN 60/20	229	186	13.2	12.6
RNN 90/10	202	173	12.8	12.2
RNN 250/5	173	155	12.3	11.7
RNN 250/2	176	156	12.0	11.9
RNN 400/10	171	152	12.5	12.1
3xRNN static	151	143	11.6	11.3
3xRNN dynamic	128	121	11.3	11.1



# Outline

---

- ▶ N-gram Models
  - ▶ Sparsity: Smoothing, Backoff
  - ▶ Perplexity Measure
- ▶ Exploiting Word Similarity
  - ▶ Brown Clustering
  - ▶ HMMs (briefly)
  - ▶ Neural Network Language Models
- ▶ **Future Directions**



---

---

