



LM EVALUATION AND NEURAL LMS

Sandeep Soni

02/27/2024

LANGUAGE MODEL

- Language modeling is a task of estimating a probabilistic model over words
- If x is a sequence, we're trying to find $P(x)$

LANGUAGE MODEL



$P(\text{"to be or not to be"})$



$P(\text{"or not to be to be"})$

- LM can answer:
- $P(x)$
- Is $P(x) > P(y)$
- Is x fluent?

N-GRAM LANGUAGE MODELS

$$P(x) = \prod_i P(x_i)$$

bigram

trigram

unigram

$$P(x) = \prod_i P(x_i | x_{i-1})$$

$$P(x) = \prod_i P(x_i | x_{i-2}, x_{i-1})$$

N-GRAM LANGUAGE MODELS

$$P(x) = \prod_i P(x_i)$$

unigram

$$P(x) = \prod_i P(x_i | x_{i-1})$$

bigram

$$P(x) = \prod_i P(x_i | x_{i-2}, x_{i-1})$$

trigram

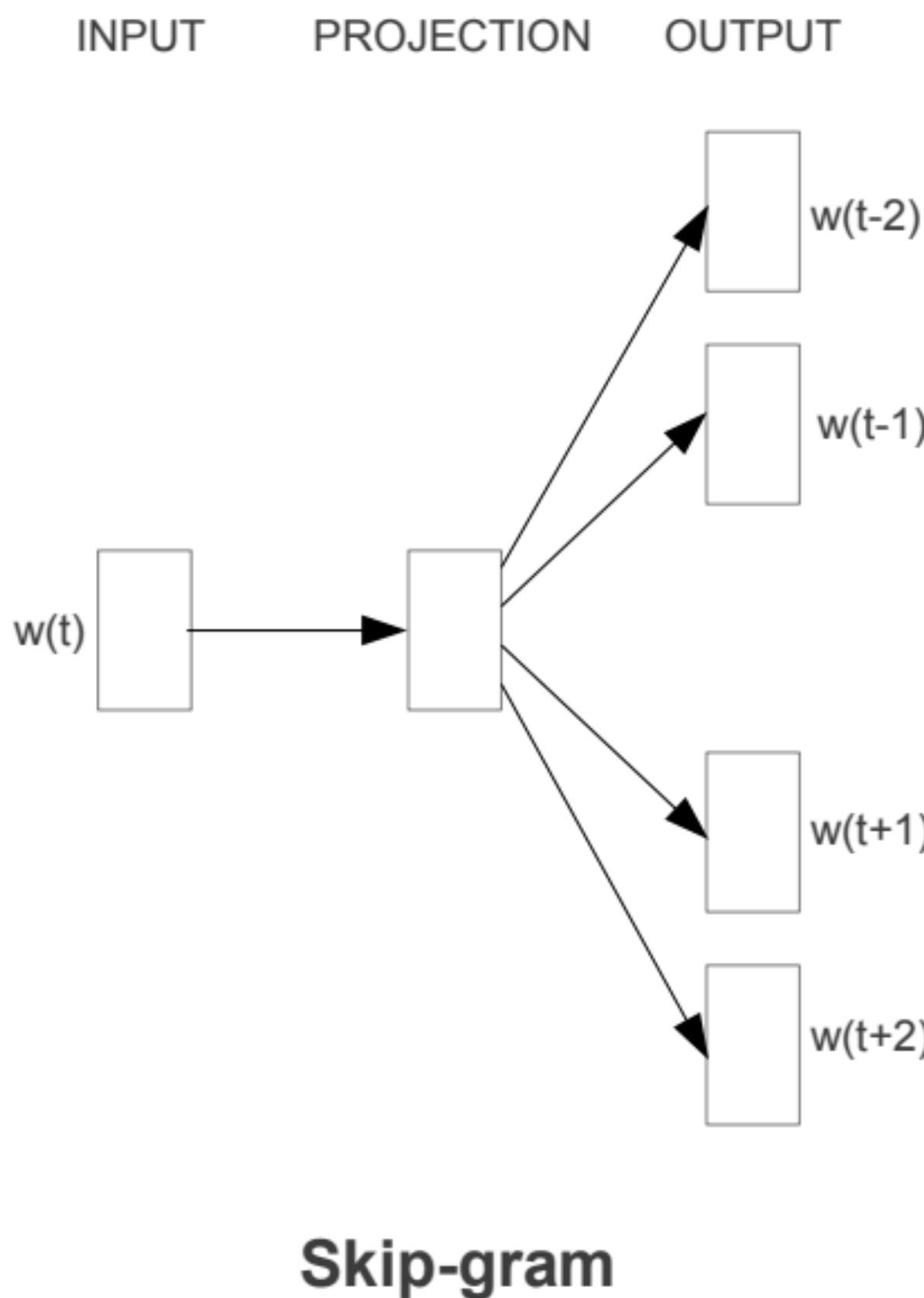
EVALUATION

- A good language model
 - assigns high probability to unseen plausible sequences
 - should be less perplexed or surprised
- Perplexity = inverse probability of test data, averaged over words

PERPLEXITY

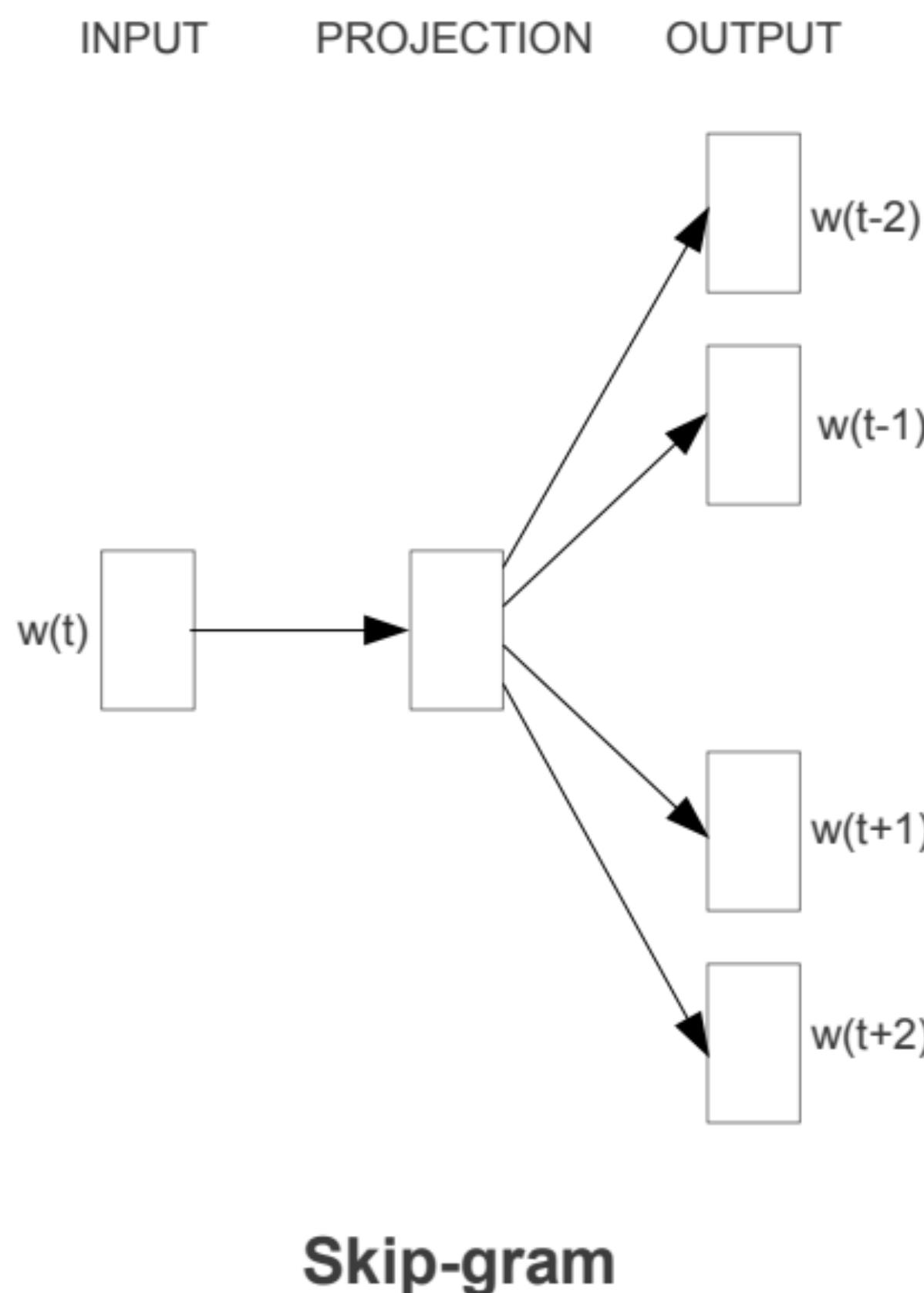
- Perplexity = $2^{-\frac{l(w)}{M}}$, where M is the total number of unseen tokens
- $$l(w) = \sum_{m=1}^M \log P(w_m | w_{m-1}, \dots, w_1)$$
- You'll calculate $l(w)$ based on the LM of your choice
- Smaller perplexity is better!

SKIPGRAM



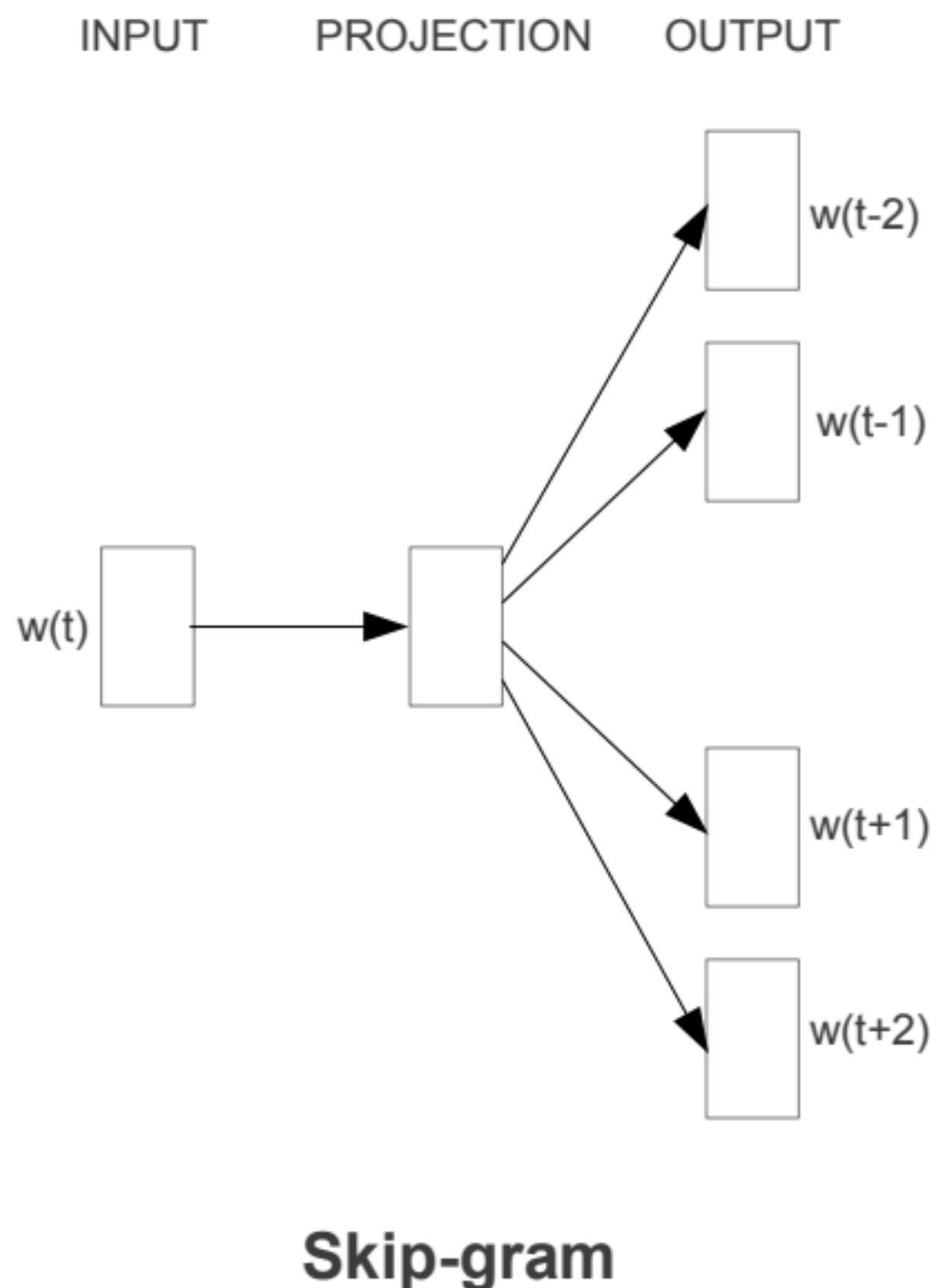
- In one of the most sensational comebacks in the champions league Liverpool defeated Barcelona by mounting a **strong** performance in the second half

SKIPGRAM



- In one of the most sensational comebacks in the champions league Liverpool defeated Barcelona by mounting a **strong** _____ in the second half

SKIPGRAM



- In one of the most sensational comebacks in the champions league Liverpool defeated Barcelona by mounting a **strong** _____ in the second half

Word2vec learned word vectors using the skipgram language model

IN CLASS

- lm exploration

No Country for Old Members: User Lifecycle and Linguistic Change in Online Communities

Cristian Danescu-Niculescu-Mizil
Stanford University
Max Planck Institute SWS
cristiand@cs.stanford.edu

Robert West
Stanford University
west@cs.stanford.edu

Dan Jurafsky
Stanford University
jurafsky@stanford.edu

Jure Leskovec
Stanford University
jure@cs.stanford.edu

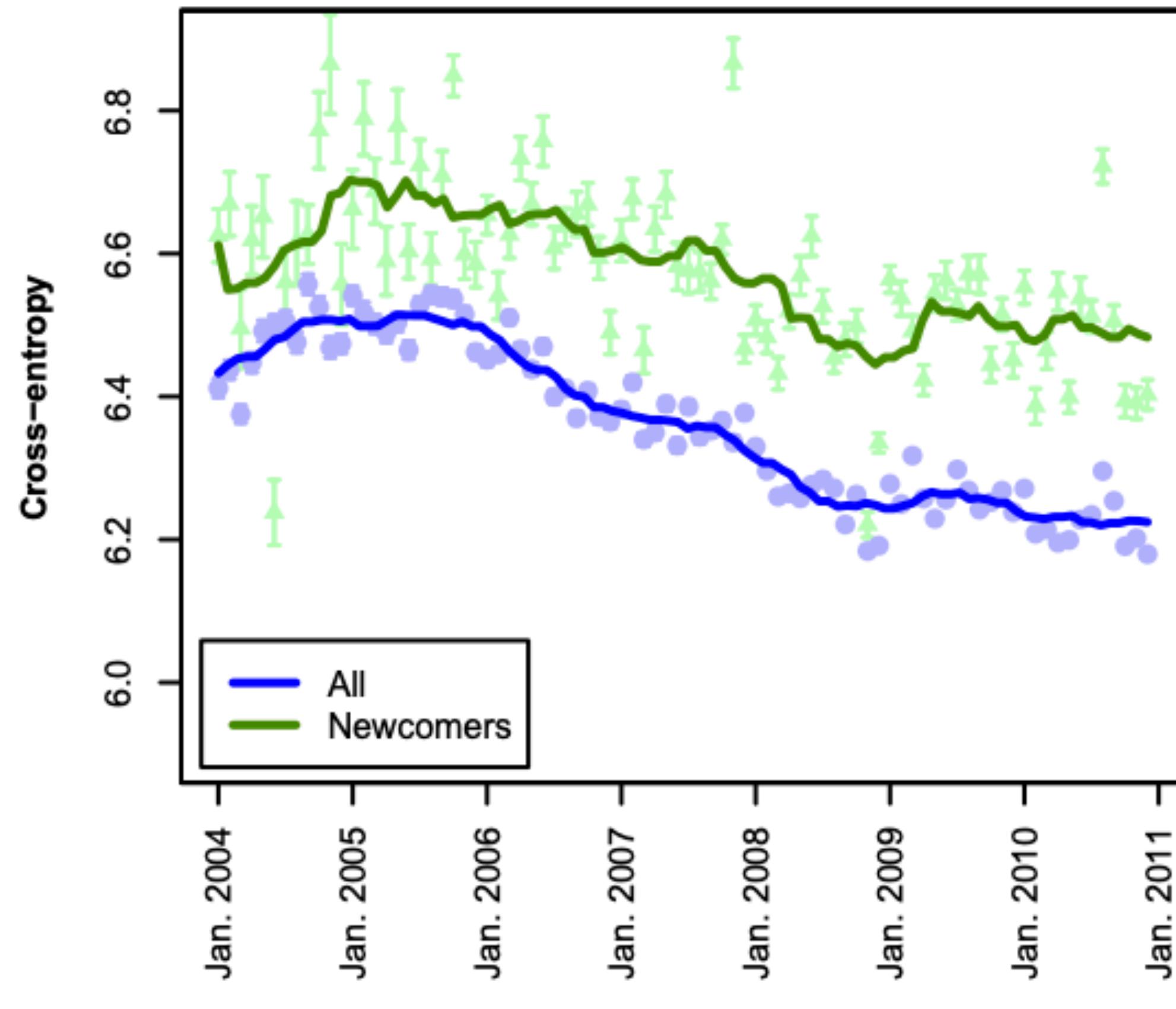
Christopher Potts
Stanford University
cgpotts@stanford.edu

Who contributes to changing linguistic norms?

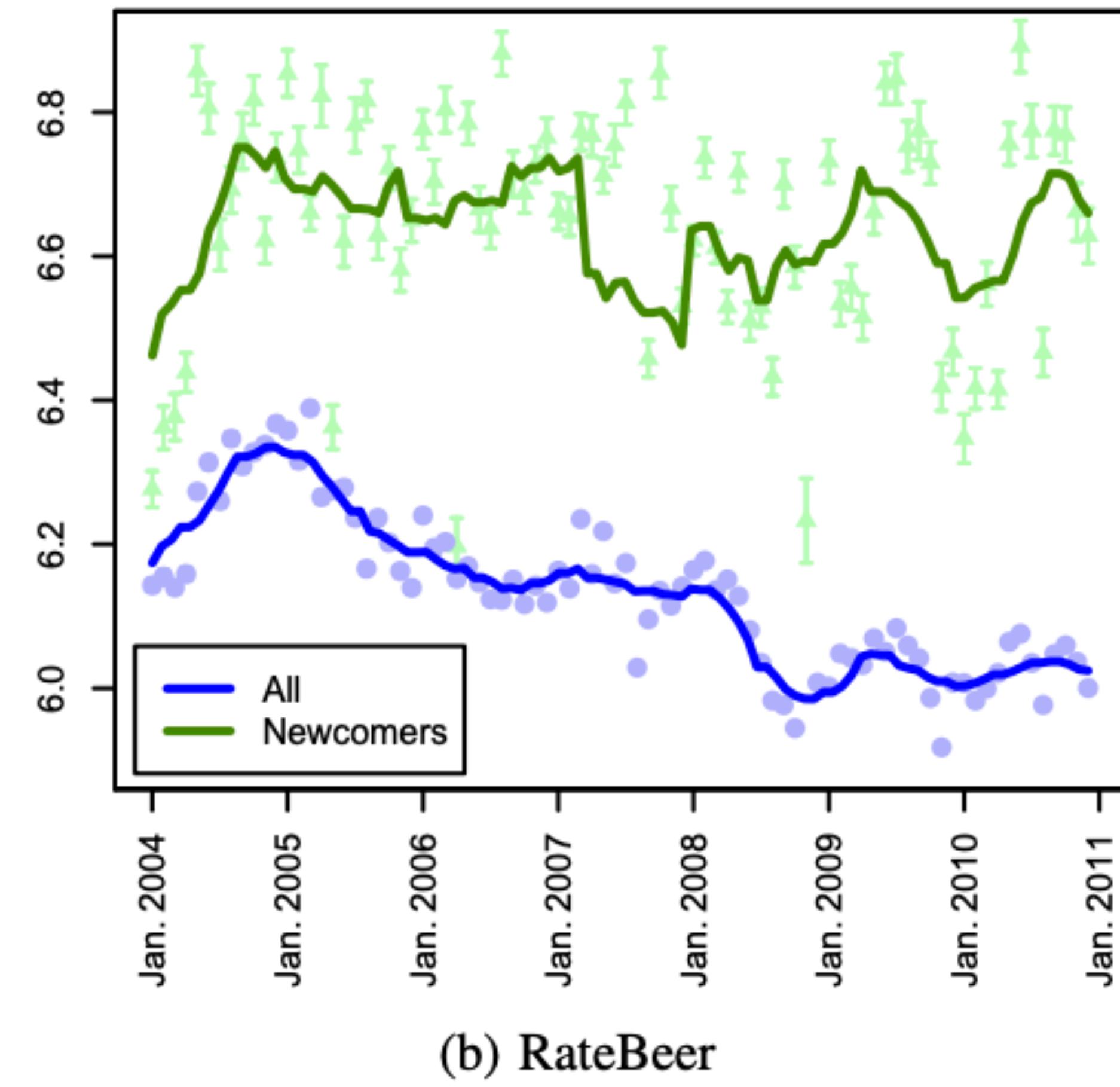
SNAPSHOT LANGUAGE MODELS

- Create a bigram language model SLM_m for every month m ; this is called a snapshot language model
- Q. How surprising is a document with respect to any month?
- Calculate cross-entropy $H(p, SLM_m) = -\frac{1}{N} \sum_i \log P_{SLM_m}(b_i)$, where

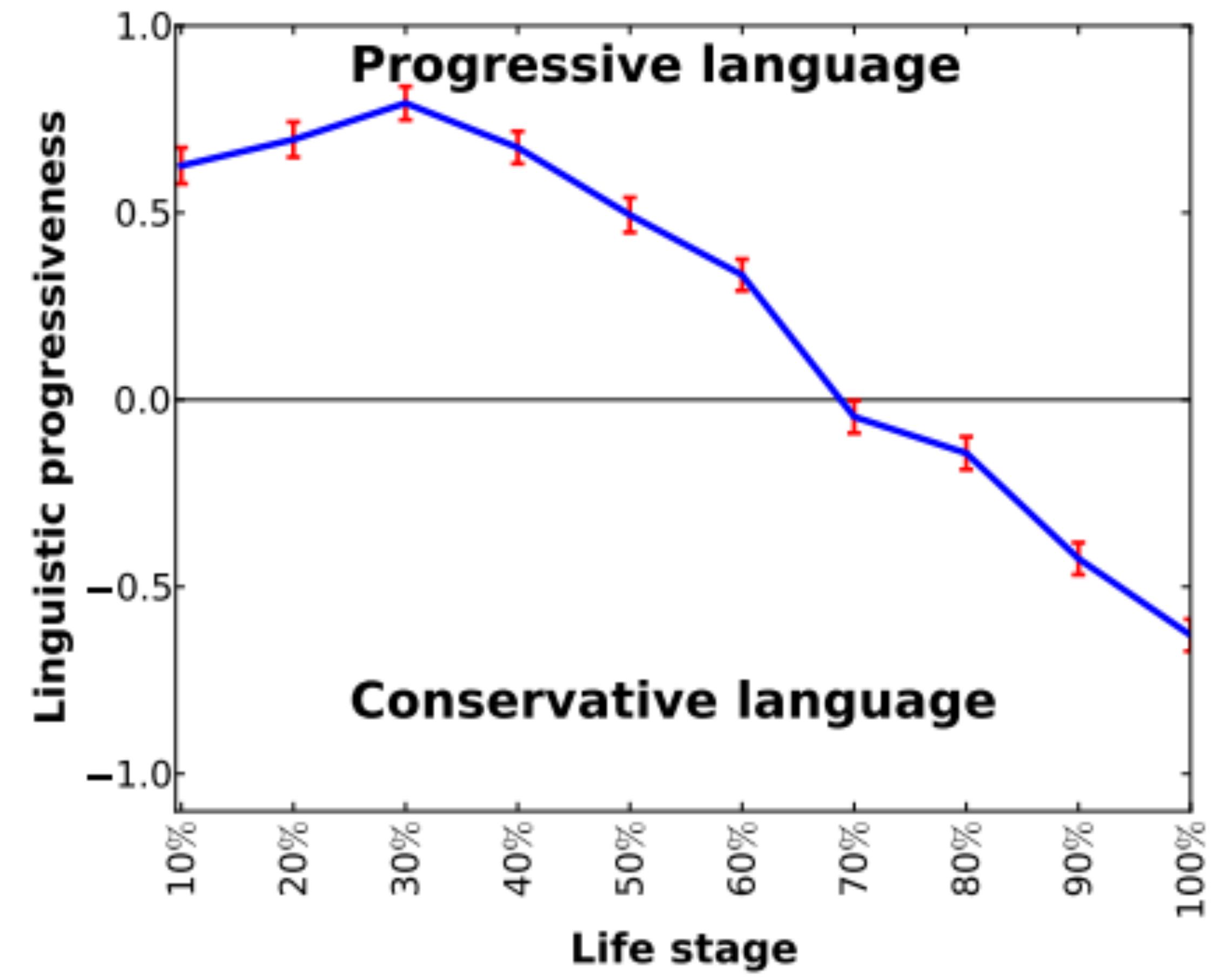
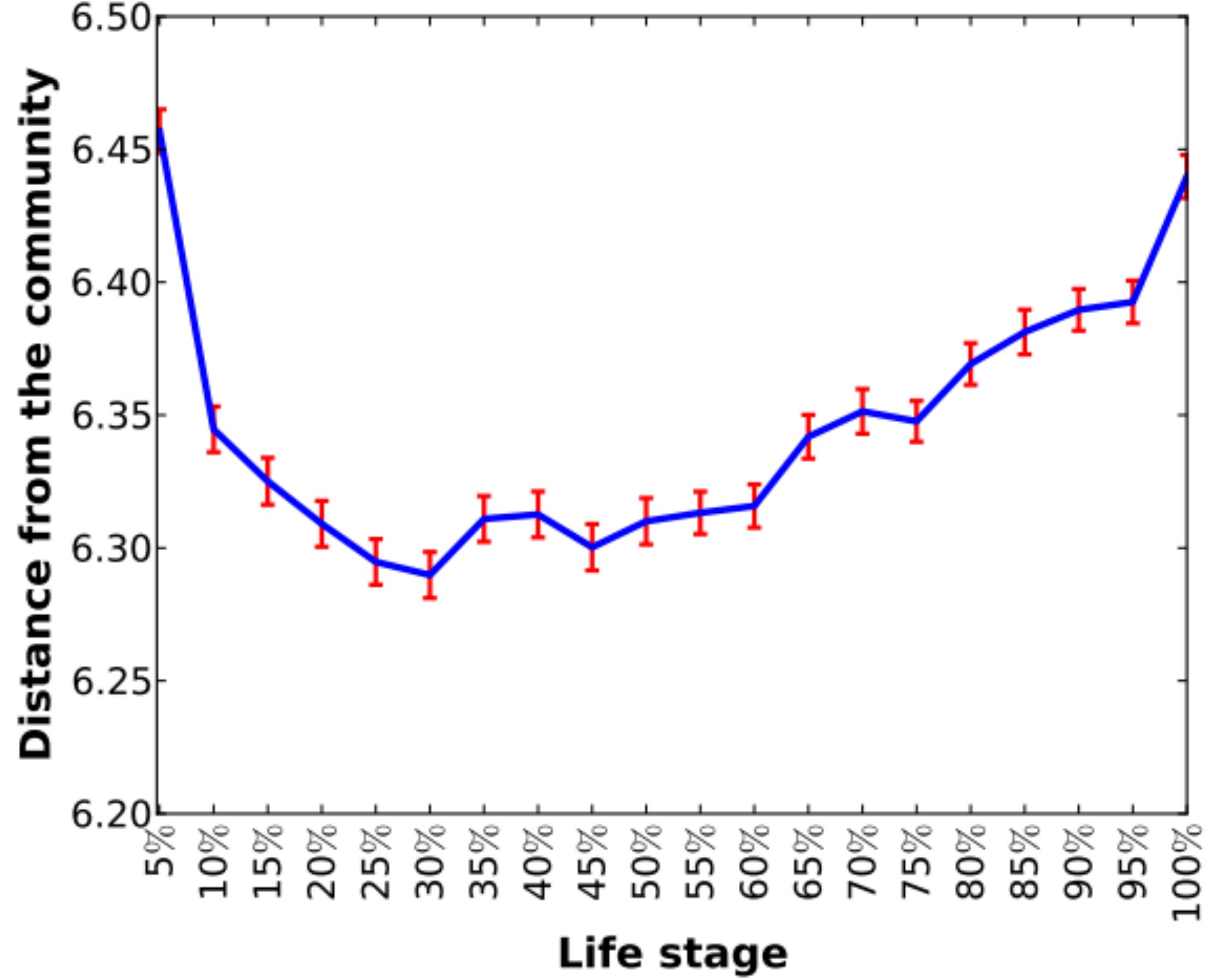
b_1, b_2, \dots, b_N are bigrams from the post



(a) BeerAdvocate

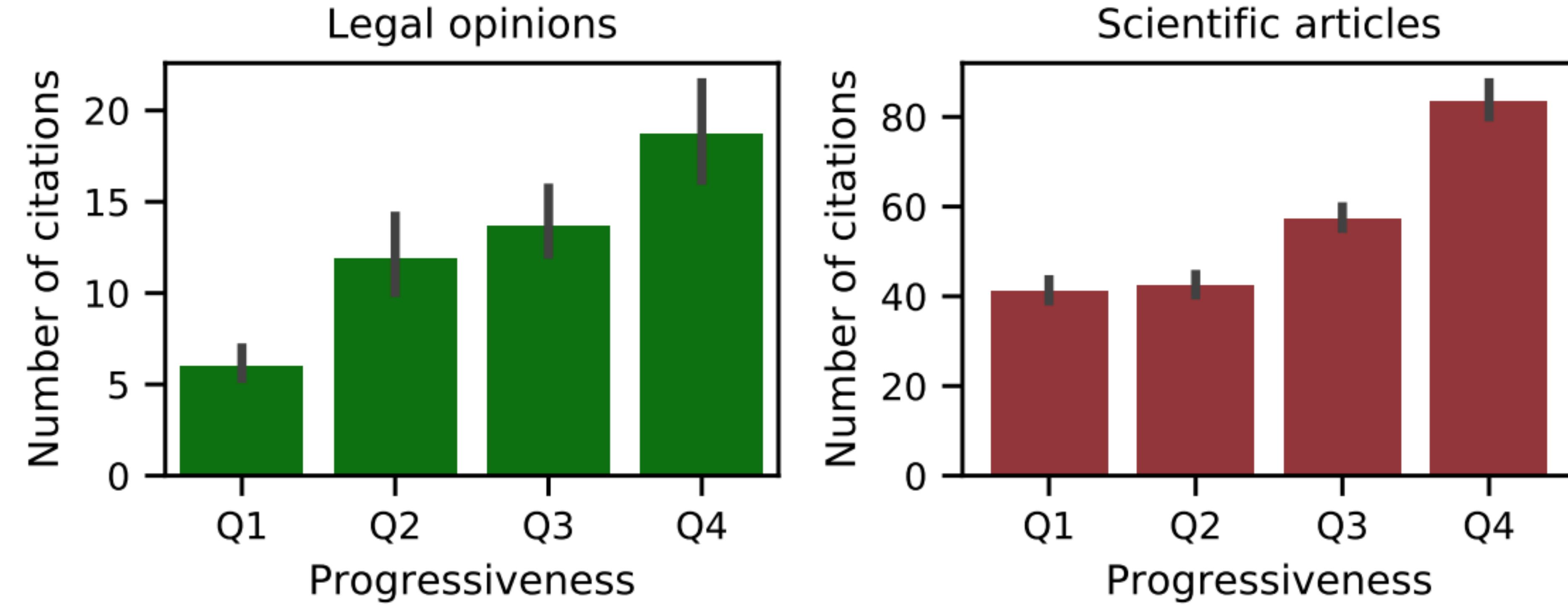


(b) RateBeer



- Users conform to the community's language initially but then stop adapting after some point

- Users are more innovative and trend-setting with their language use initially but then stabilize to rely increasingly more on past language



- A legal opinion or a scientific article is progressive if it uses words with meanings from the future than the past
- Learn two skipgram language models, one trained on future data and one trained on the past; then compare word senses with respect to both the models

Can we do better than N-gram language models?

GENERATIVE VS DISCRIMINATIVE

	Naive Bayes	Logistic Regression
Type of classifier	Generative	Discriminative
Model	$P(x,y)$	$P(y \mid x)$
Objective	Generate the data and label jointly	Predict the label from the data

Don't count, predict! A systematic comparison of context-counting vs. context-predicting semantic vectors

Marco Baroni and Georgiana Dinu and Germán Kruszewski

Center for Mind/Brain Sciences (University of Trento, Italy)

(marco.baroni|georgiana.dinu|german.kruszewski)@unitn.it

Abstract

Context-predicting models (more commonly known as embeddings or neural language models) are the new kids on the distributional semantics block. Despite the buzz surrounding these models, the literature is still lacking a systematic comparison of the predictive models with classic, count-vector-based distributional semantic approaches. In this paper, we perform such an extensive evaluation, on a wide range of lexical semantics tasks and across many parameter settings. The results, to our own surprise, show that the buzz is fully justified, as the context-predicting models obtain a thorough and resounding victory against their count-based counterparts.

1 Introduction

A long tradition in computational linguistics has shown that contextual information provides a good approximation to word meaning, since semanti-

optimization process is generally unsupervised, and based on independent considerations (for example, context reweighting is often justified by information-theoretic considerations, dimensionality reduction optimizes the amount of preserved variance, etc.). Occasionally, some kind of indirect supervision is used: Several parameter settings are tried, and the best setting is chosen based on performance on a semantic task that has been selected for tuning.

The last few years have seen the development of a new generation of DSMs that frame the vector estimation problem directly as a supervised task, where the weights in a word vector are set to maximize the probability of the contexts in which the word is observed in the corpus (Bengio et al., 2003; Collobert and Weston, 2008; Collobert et al., 2011; Huang et al., 2012; Mikolov et al., 2013a; Turian et al., 2010). The traditional construction of context vectors is turned on its head: Instead of first collecting context vectors and then reweighting these vectors based on various criteria, the vector weights are directly set to optimally predict the contexts in which the corresponding

Learn
representations
from data by
predicting parts
of the data,
instead of
counting

LANGUAGE MODELING



- Instead of modeling $P(x)$, why not model $P(w|c)$?
- Rather than directly estimating the word probabilities from relative frequencies, this hints at language modeling as a learning task

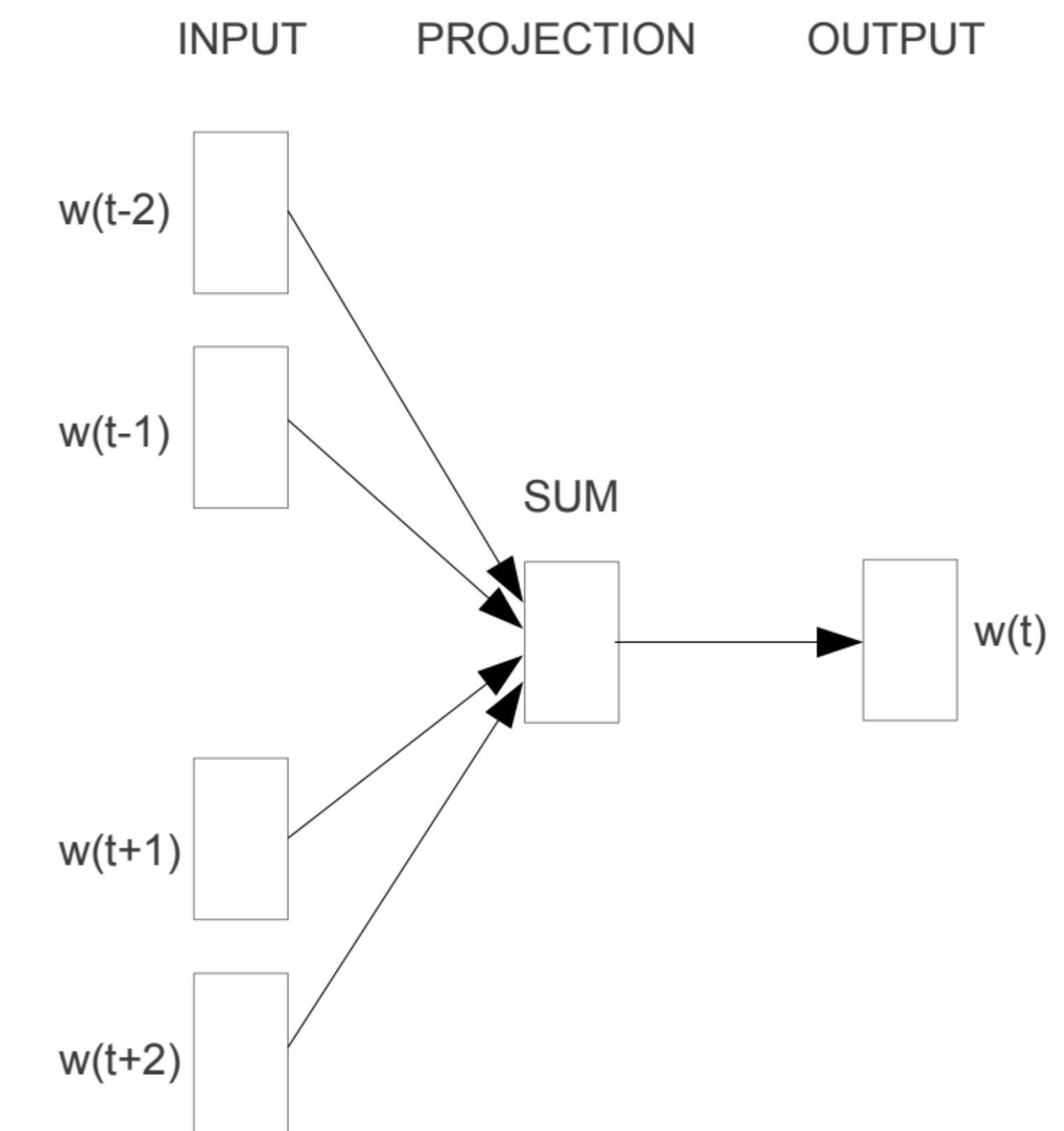
LANGUAGE MODELING



- Reparametrize the probability $P(w|c)$ to depend on dense vectors
- $$P(w|c) = \frac{\exp(\beta_w v_c)}{\sum_{w'} \exp(\beta_{w'} v_c)}$$
, where β_w is a vector representation of w and v_c is the vector representation of the context

WORD2VEC (CBoW)

- $$P(w|c) = \frac{\exp(\beta_w v_c)}{\sum_{w'} \exp(\beta_{w'} v_c)}$$
- In CBoW model of word2vec, w is a word and c are words on the left and right of w
- v_c was calculated as a sum of output vectors

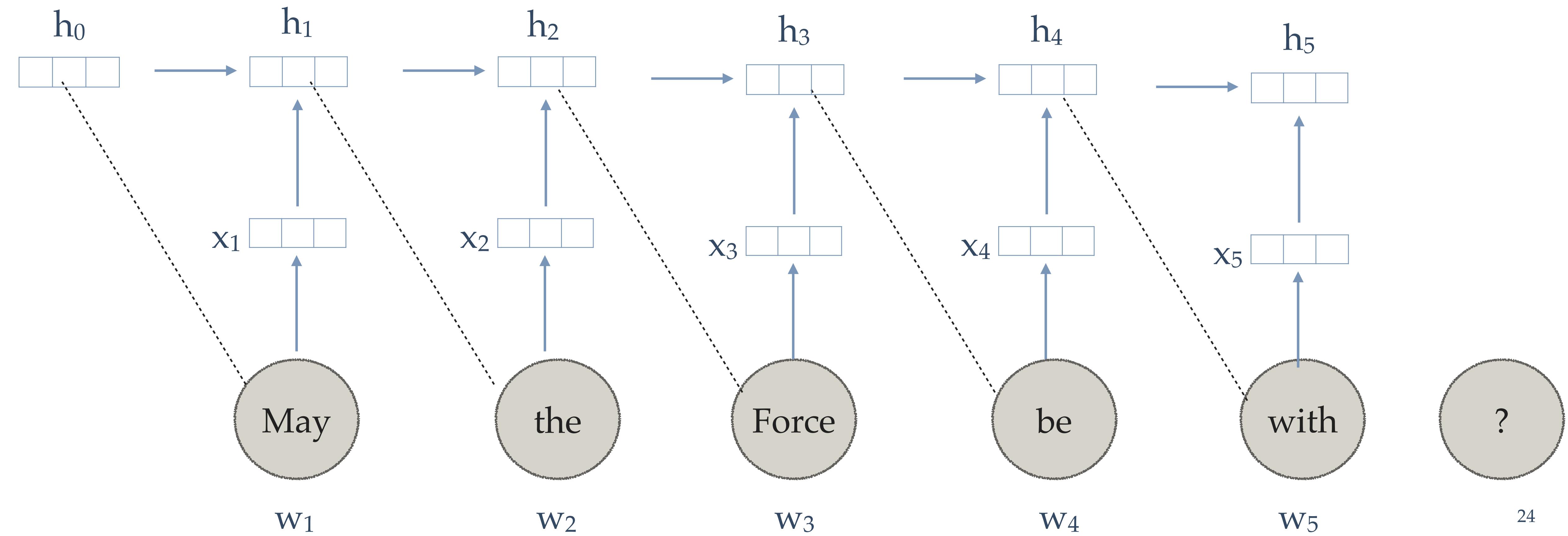
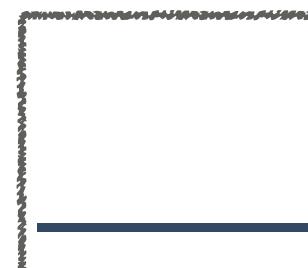


What can be better ways of coming up with a vector representation of the context?

Context

Word

May the Force be with



RECURRENT NEURAL NETWORK LM

At every position m:

$$x_m = \text{Lookup}(\phi, w_m)$$

$$h_m = \text{RNN}(x_m, h_{m-1})$$

$$h_m = g(\Theta h_{m-1} + x_m)$$

Elman unit

$$P(w_{m+1} | w_1, w_2, \dots, w_m) = \text{softmax}(\beta_{w_m}, \mathbf{h}_m)$$

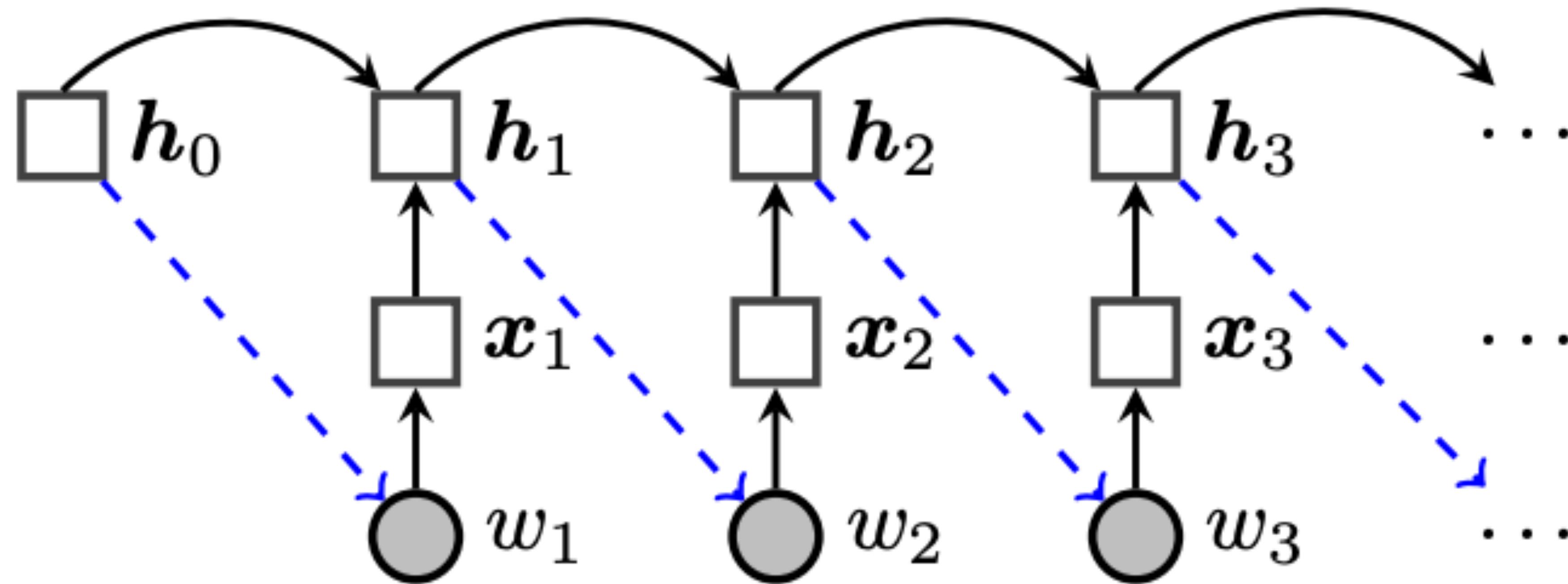
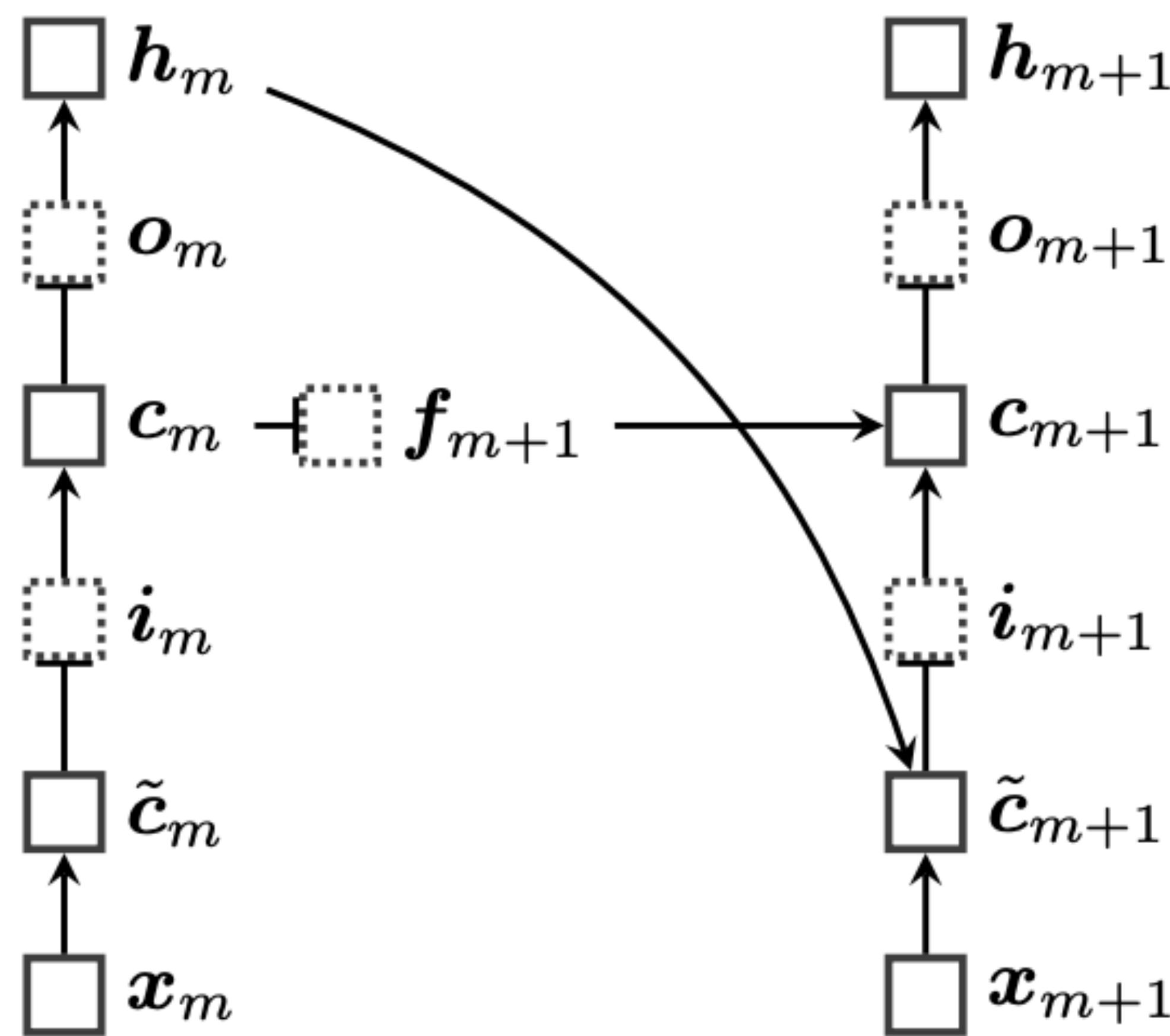


Figure taken from Eisenstein 2018

LONG SHORT-TERM MEMORIES (LSTM)



- Transform x to h by passing x through gating units
- Preserves information propagation over long distances and downweights unimportant contexts in the past

SUMMARY

- Language Modeling is a foundational task in text processing
- Count based language models are easy to interpret but not very powerful for longer sequences
- Neural LMs (RNNs, LSTMs) are extremely powerful
- Can we improve any further?