



Chapter 3:

Basics of Language Modelling



Motivation

Language Models are used in

- Speech Recognition
- Machine Translation
- Natural Language Generation
- Query completion
- Caption Generation

Quality of LMS:

need a simple evaluation metric for fast turn-around times in experiments



Test your Language Model



What's in your hometown newspaper ???



Test your Language Model



What's in your hometown newspaper **today**



Test your Language Model



It's raining cats and ???



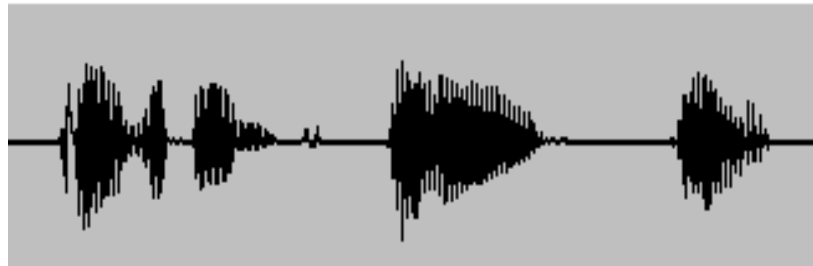
Test your Language Model



It's raining cats and **dogs**



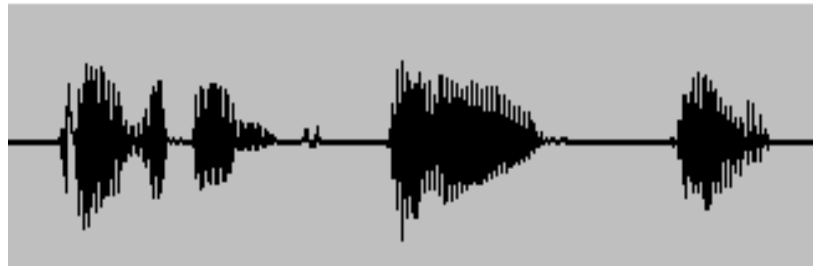
Test your Language Model



President Bill ???



Test your Language Model



President Bill Gates



Definition Language Model

A language model either:

- Predicts the next word w given a sequence of predecessor words h (history)

$$P(w_i | h_i) = P(w_i | w_1, \dots, w_{i-1}) \text{ with } i : \text{position in text}$$

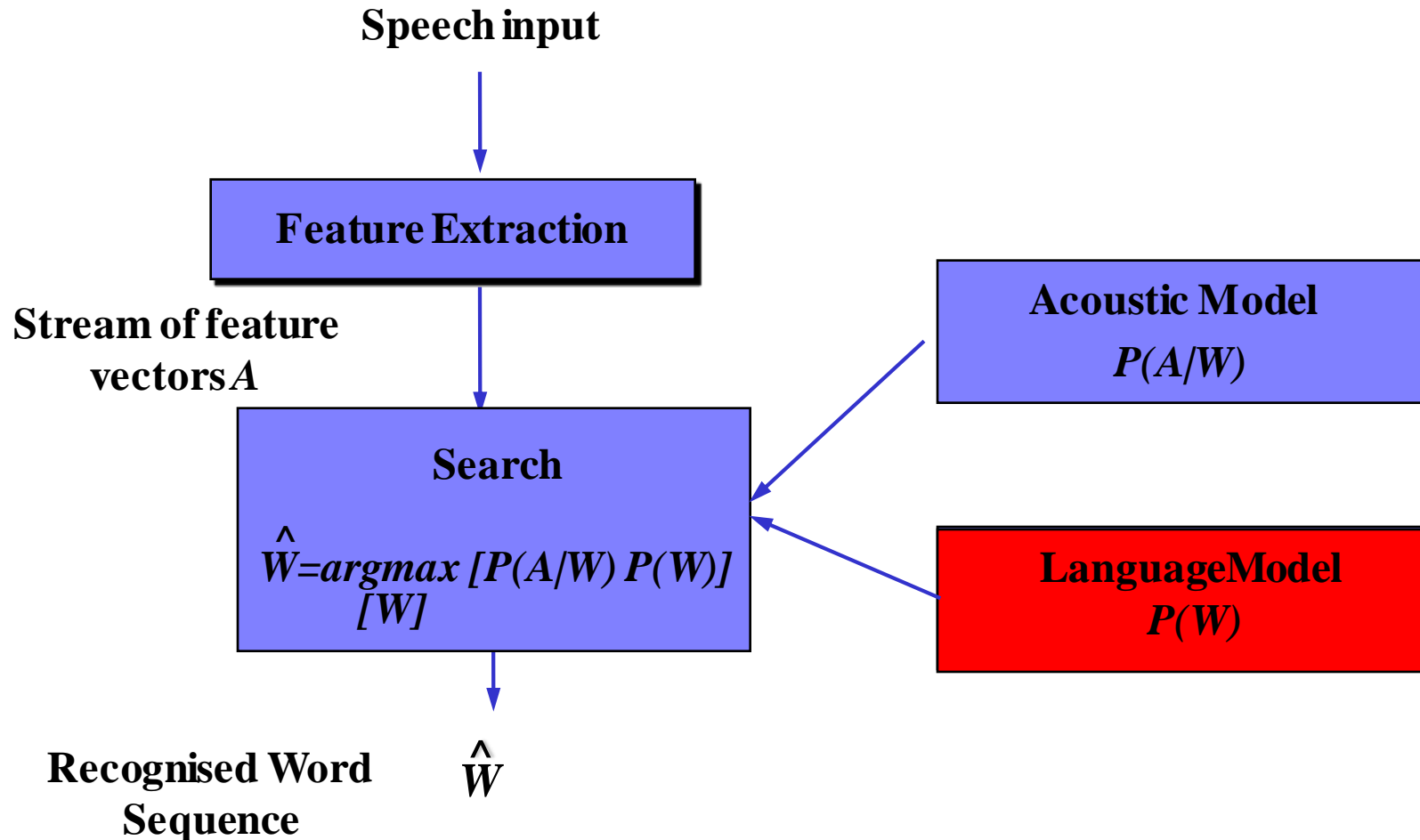
or

- Predicts the probability of a sequence of words

$$P(w_1, w_2, \dots, w_{N-1}, w_N) = P(W)$$

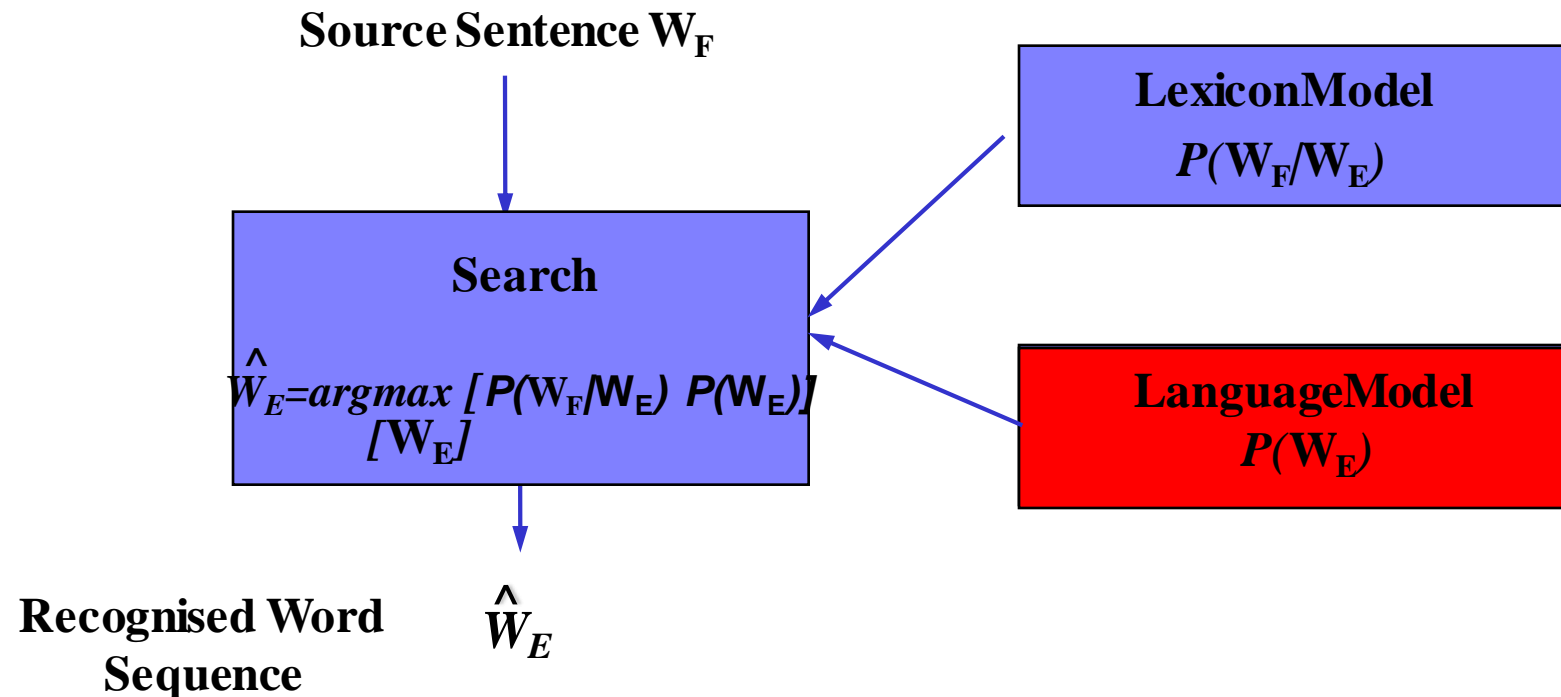


Language Model in Speech Recognition (ASR)

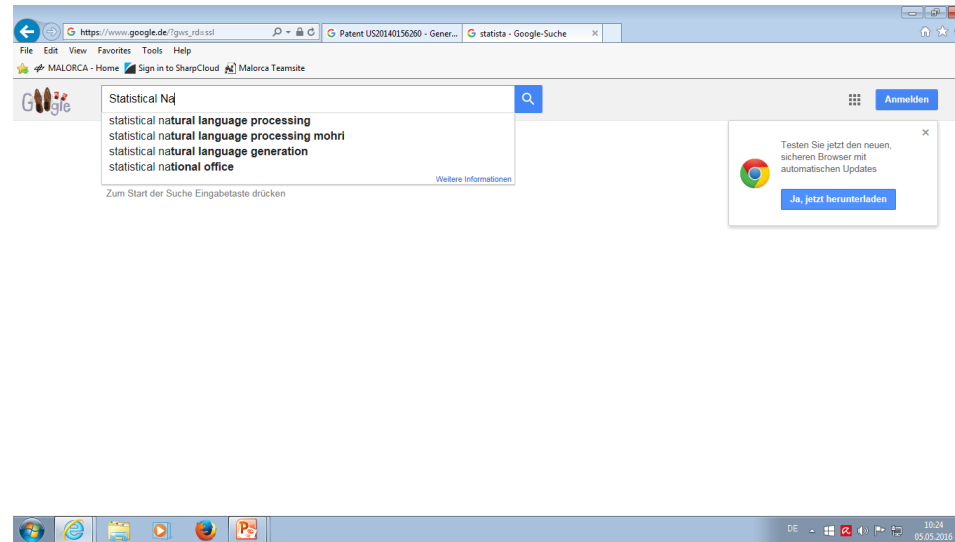




Language Model in Machine Translation (MT)



Sentence Completion



- Directly use of $P(w|h)$
- Needs good search algorithms



The Markov Assumption: truncate the history

Cutting history to M-1 words:

$$P(w_i | w_1, w_2, \dots, w_{i-1}) = P(w_i | w_{i-M+1}, \dots, w_{i-1})$$

Special cases:

“Zerogram”
(uniform distribution)

$$P(w_i | w_1, w_2, \dots, w_{i-1}) \approx \frac{1}{W}$$

Unigram

$$P(w_i | w_1, w_2, \dots, w_{i-1}) \approx P(w_i)$$

Bigram

$$P(w_i | w_1, w_2, \dots, w_{i-1}) \approx P(w_i | w_{i-1})$$

Trigram

$$P(w_i | w_1, w_2, \dots, w_{i-1}) \approx P(w_i | w_{i-2}, w_{i-1})$$

....

How would
you measure
the quality of
an LM?



Definition of Perplexity

Let w_1, w_2, \dots, w_N be an independent test corpus not used during training.

Perplexity

$$PP = P(w_1, w_2, \dots, w_N)^{-1/N}$$

Interpretation:

Normalized probability of test corpus



Example: Zerogram Language Model

“Zerogram”
(uniform distribution) $P(w_i | w_1, w_2, \dots, w_{i-1}) \approx \frac{1}{W}$

$$PP = P(w_1, w_2, \dots, w_N)^{-1/N}$$

$$= \left(\prod_{i=1}^N \frac{1}{W} \right)^{-1/N}$$

$$= W$$

Interpretation: perplexity is the
“average” de-facto size of vocabulary



Alternate Formulation

Assumption:

Use an M-gram language model

History $h_i = w_{i-M+1} \dots w_{i-1}$

Idea:

collapse all identical histories



Alternate Formulation

Example:

$$M=2$$

Corpus:

“to be or not to be”

$$h_2 = \text{“to”} \quad w_2 = \text{“be”}$$

$$h_6 = \text{“to”} \quad w_6 = \text{“be”}$$

↪ calculate $P(\text{“be”}|\text{“to”})$ only once and scale it by 2



Alternate Formulation

$$PP = P(w_1, w_2, \dots, w_N)^{-1/N}$$

Use Bayes decomposition

$$= \left(\prod_{i=1}^N P(w_i | h_i) \right)^{-1/N}$$

Try to get rid of product

$$= \exp \left(\log \left(\prod_{i=1}^N P(w_i | h_i) \right)^{-1/N} \right)$$

$$= \exp \left(-\frac{1}{N} \log \left(\prod_{i=1}^N P(w_i | h_i) \right) \right)$$



Alternate Formulation

$$= \exp \left(-\frac{1}{N} \log \left(\prod_{i=1}^N P(w_i | h_i) \right) \right)$$

$$= \exp \left(-\frac{1}{N} \sum_{i=1}^N \log(P(w_i | h_i)) \right)$$

$$= \exp \left(-\frac{1}{N} \sum_{w,h} N(w,h) \log(P(w|h)) \right)$$

$N(w,h)$: absolute frequency
of sequence h,w on test corpus

$$= \exp \left(-\sum_{w,h} f(w,h) \log(P(w|h)) \right)$$

$f(w,h)$ relative frequency
of sequence h,w on test corpus



Alternate Formulation: Final Result



$$PP = P(w_1 \dots w_N)^{-1/N}$$
$$= \exp \left(- \sum_{w,h} f(w,h) \log(P(w|h)) \right)$$

Perplexity can also be calculated using conditional probabilities trained in the training corpus and relative frequencies from the test corpus.



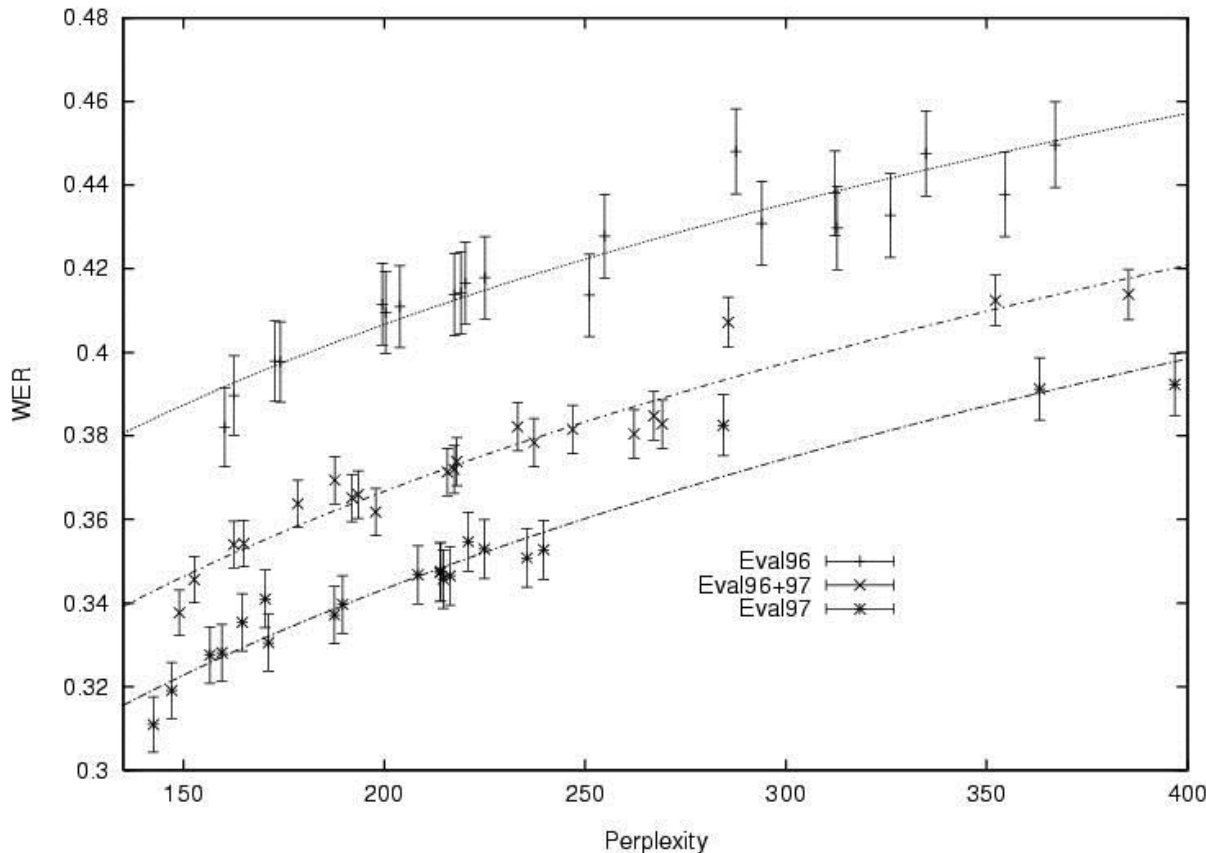
Alternate Formulation: Zero-gram Example

“Zero-gram”
(uniform distribution) $P(w_i | w_1, w_2, \dots, w_{i-1}) \approx \frac{1}{W}$

$$\begin{aligned} PP &= \exp \left(- \sum_{w,h} f(w,h) \log(P(w|h)) \right) = \exp \left(- \sum_{w,h} f(w,h) \log \left(\frac{1}{W} \right) \right) \\ &= \exp \left(- \log \left(\frac{1}{W} \right) \sum_{w,h} f(w,h) \right) = \exp (\log(W) * 1) = W \end{aligned}$$

Identical result

Perplexity and Error Rate



Perplexity is
correlated to
word error rate

Power law relation



Quality Measure: Mean Rank

- Definition:
 - Let w follow h in the test text
 - Sort all words after a given history according to $p(w|h)$
 - Determine position of correct word w
 - Average over all events in the test text



Alternative: Average Rank

1 The are to know the issues necessary
2 This will the have this problems data
3 One the would understand do these above
4 Two do also do get the problems any other
5 A do do the use a problem them time
6 Three need do the use a problem them people
7 Please do do the use a problem them operators
8 In do do the use a problem them
9 We do do the use a problem them tools
.
.
.
93 request
94 respond
95 supply
96 write
97 me
98 resolve
.
.
.
1636
1637
1638
1639
1640
1641

Task:
guess the next word

Metric:
Measure how many
attempts it takes

1 role and the next be meeting of
2 thing from
3 that in
4 to to
5 contact are
6 parts with
7 point were
8 for requiring
9 issues still
.
.
.
61 being
62 during
63 I
64 involved
65 would
66 within



Quality Measure: Mean Rank

Measure	Correlation with ASR error rate
Perplexity	0.955
Mean Rank	0.957

Mean rank equally good in predicting ASR performance



Summary

- Language models predict the next word
- Applications:
 - Speech Recognition
 - Machine Translation
 - Natural Language Generation
 - Information retrieval
 - Sentence Completion
- Perplexity:
 - Measures quality of language model