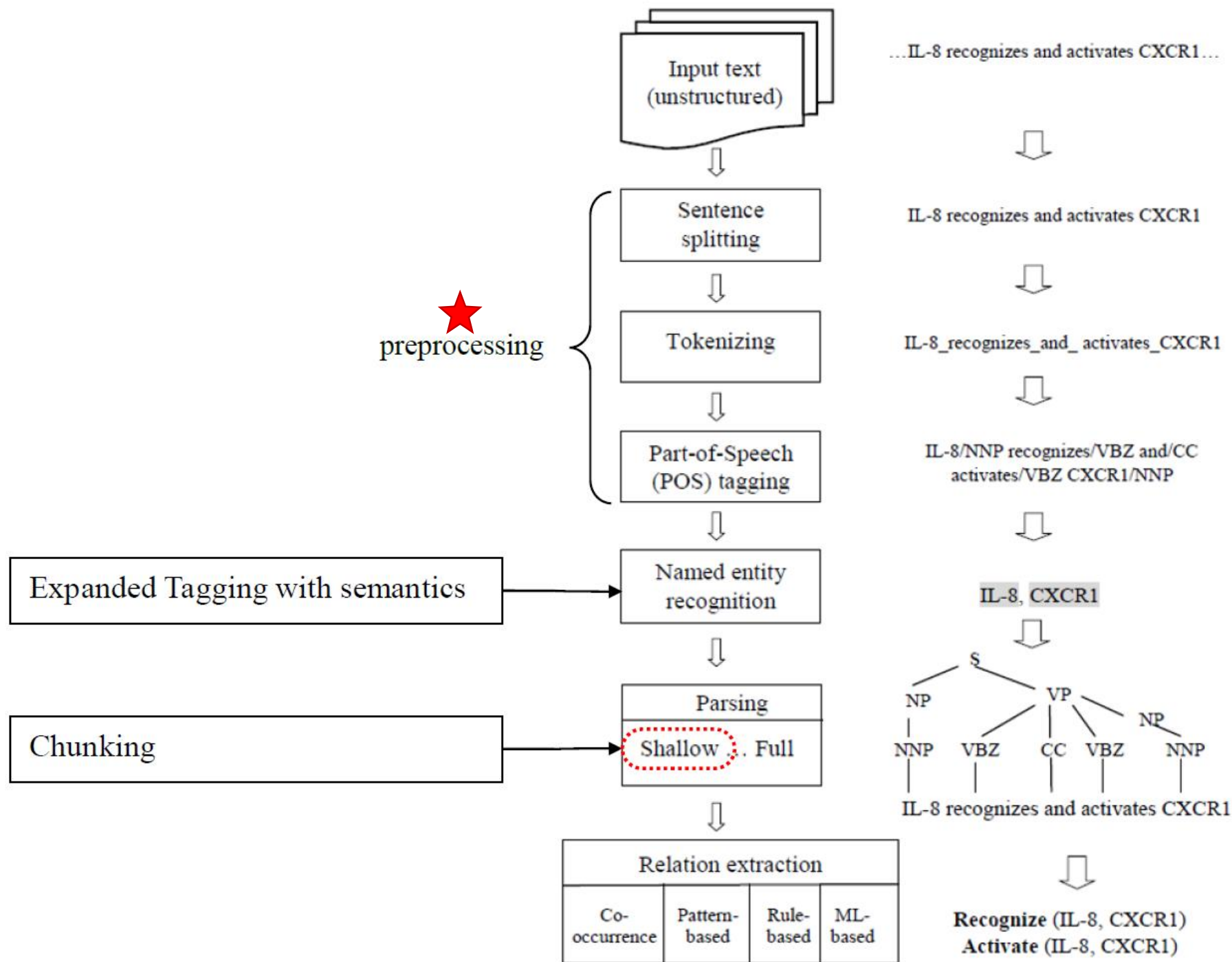


Preprocessing(2)



*n*-gram

# $n$ -gram

**WHAT**     **Contiguous sequence** of  $n$  items from a given sample of text or speech  
*phonemes, syllables, letters, words* or *base pairs* according to the application

**WHY**     **Predicts** a next letter

Given a sequence of letters, what is the likelihood of the next letter?



$x_{i-(n-1)}, \dots, x_{i-1}$



$x_i$

$P(x_i | x_{i-(n-1)}, \dots, x_{i-1})$

# WPM

## Word Piece Model

하나의 단어를 내부 단어(Subword Unit)들로 분리하는 단어 분리 모델

*Google's Neural Machine Translation System:  
Bridging the Gap between Human and Machine Translation  
Wot et al, 2016*

System	Combined test set
Bing Translate	23.48
Google Translate	22.67
<b>Word Model</b>	18.19
<b>Word-piece Model</b>	26.58

# BPE

## Byte Pair Encoding (Digram Coding)

Simple form of data compression

The **most common pair** of consecutive bytes of data is replaced with a byte

*Philip Gage, 1994*

*Neural Machine Translation of Rare Words with Subword Units*

*Sennrich et al, 2015*

# example

- learning
    - word:freq : {low:5, lowest:2, newer:6, wider:3}
    - marge & count
      1. 'r' '</w>' : 9 → marge'r</w>'
      2. 'e' 'r</w>' : 9 →marge'er</w>'
      3. 'l' 'o' : 7 →marge'lo'
      4. 'lo' 'w' : 7 →marge'low'
- OOV : 'lower' segmented 'low er</w>'

# Empirical Law



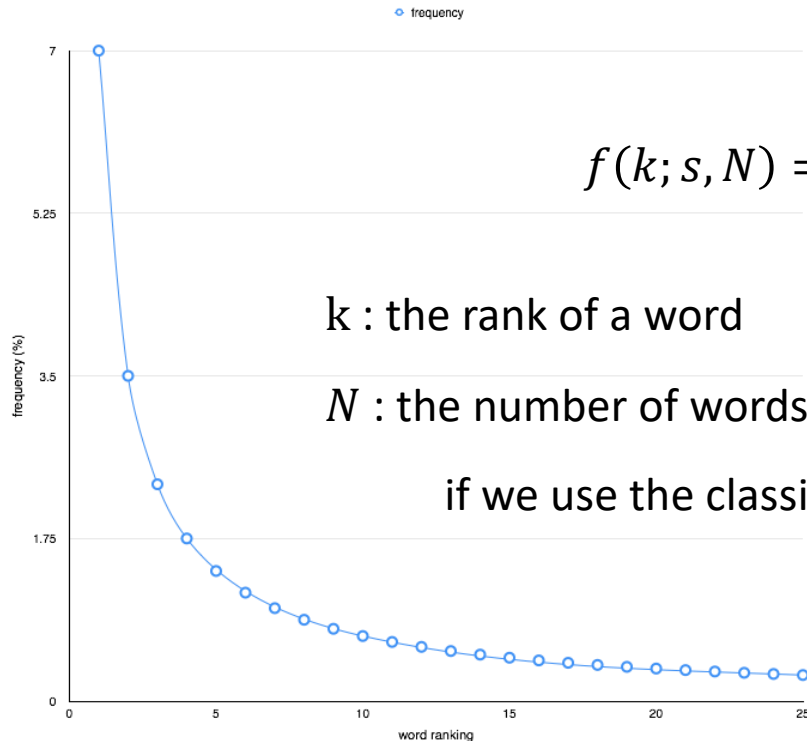
# Zipf's Law

## WHAT An empirical law

frequency of word is inversely proportional to its rank in frequency table

most frequent word will occur approximately twice

as often as the second most frequent word

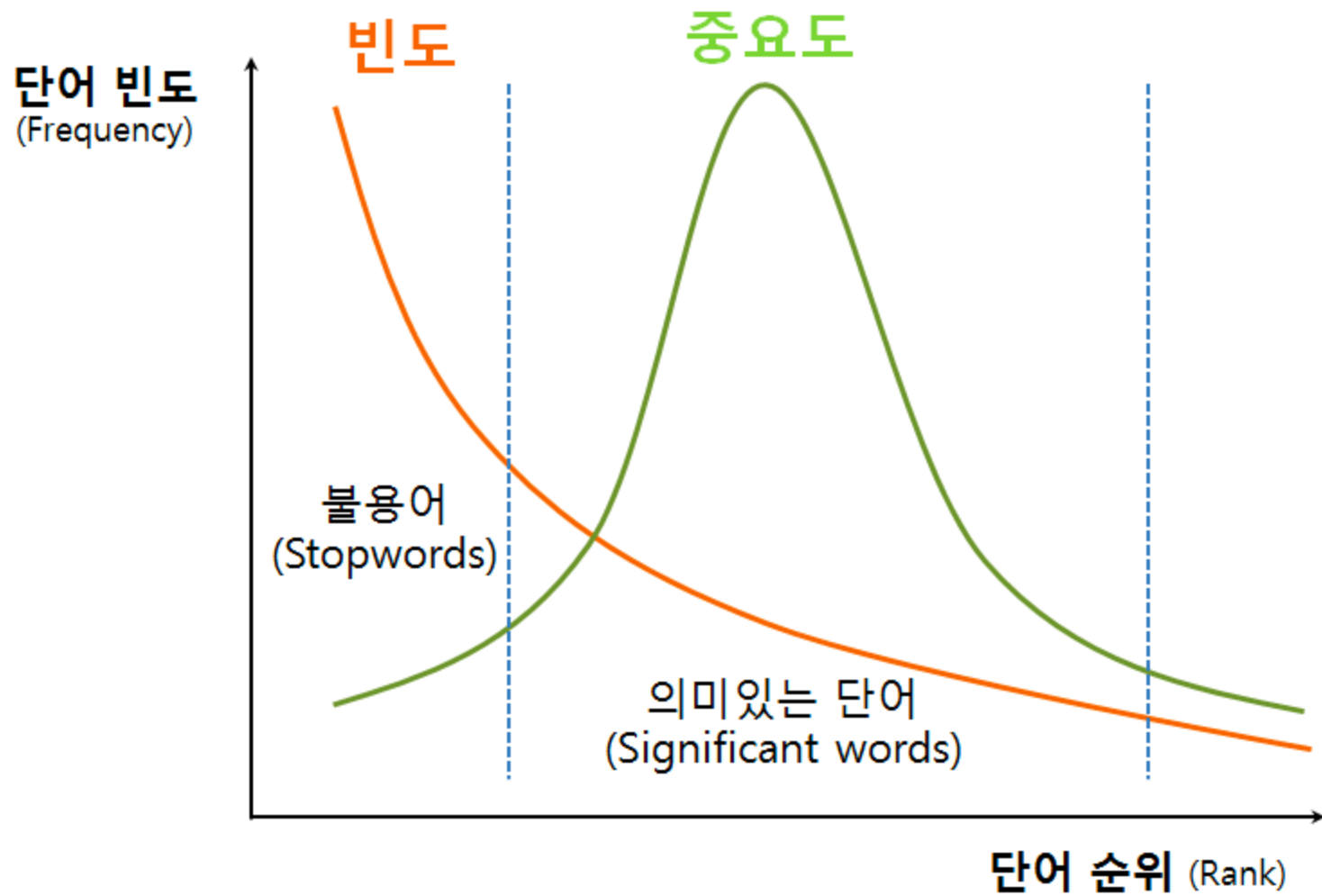


$$f(k; s, N) = \frac{1/k^s}{\sum_{n=1}^N (1/n^s)} \Leftrightarrow \frac{1}{k^s H_{n,s}}$$

$k$  : the rank of a word

$N$  : the number of words

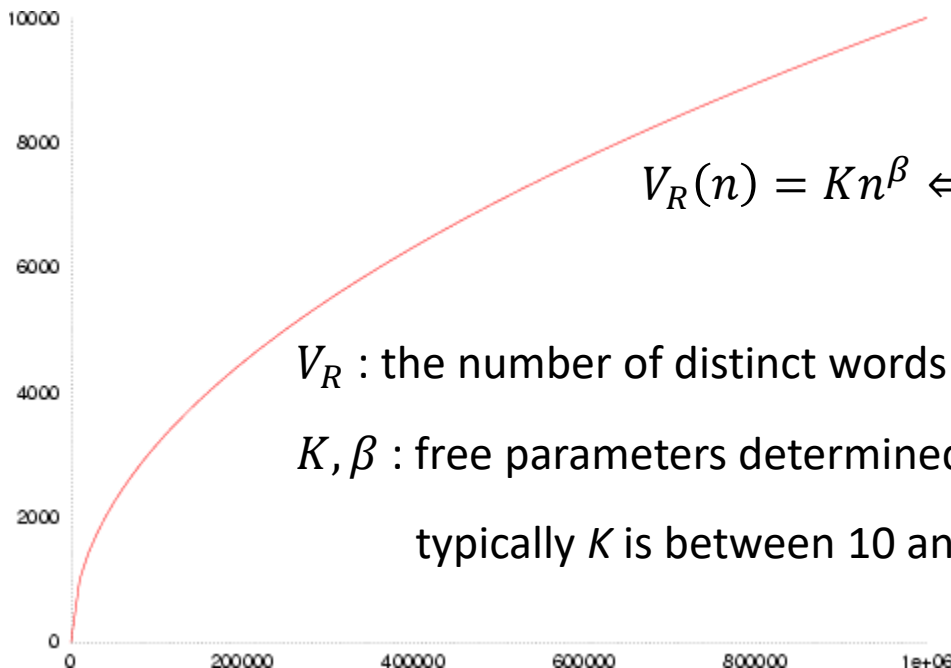
if we use the classic version of Zipf's law, the exponent  $s$  is 1



# Heaps' Law

## WHAT An empirical law

# of distinct words in a document as a function of the document length



$$V_R(n) = Kn^\beta \Leftrightarrow M = kT^b$$

$V_R$  : the number of distinct words in an instance text of size  $n$

$K, \beta$  : free parameters determined empirically

typically  $K$  is between 10 and 100, and  $\beta$  is between 0.4 and 0.6