



# Natural Language Processing

## Lecture 01

Qun Liu, Valentin Malykh  
Huawei Noah's Ark Lab



Autumn 2020  
A course delivered at MIPT, Moscow



# Content

- 1 About the course
- 2 Research questions and NLP tasks
- 3 Grammars and Automata
- 4 Text segmentation and morphology analysis
- 5 Word frequency and collocations



# Content

- 1 About the course
- 2 Research questions and NLP tasks
- 3 Grammars and Automata
- 4 Text segmentation and morphology analysis
- 5 Word frequency and collocations



# Logistics

- Instructors: Dr. Valentin Malykh
- TAs: TBA
- Time: 18.30, Tuesday
- Location: Online
- Slides: will be uploaded to the course platform before each class.



# Logistics

- Course platform is Stepik.org
- You will be invited by email you have registered with.
- The official support channel will be at OpenDataScience slack **#huawei\_nlp\_course**

# Logistics

ods.ai





# Course description

**Natural Language Processing (NLP)** is a domain of research whose objective is to analyze and understand human languages and develop technologies to enable human machine interactions with natural languages. NLP is an interdisciplinary field involving linguistics, computer sciences and artificial intelligence. The goal of this course is to provide students with comprehensive knowledge of NLP. Students will be equipped with the principles and theories of NLP, as well as various NLP technologies, including rule-based, statistical and neural network ones. After this course, students will be able to conduct NLP research and develop state-of-the-art NLP systems.



# Grading policy

- Quizzes – 30%
  - 1st Assignment – 15%
  - 2nd Assignment – 20%
  - Final project – 35%
- 
- To pass the course you need at least 80%.





# Content

- 1 About the course
- 2 Research questions and NLP tasks**
- 3 Grammars and Automata
- 4 Text segmentation and morphology analysis
- 5 Word frequency and collocations



# Natural language processing in Wikipedia

**Natural language processing (NLP)** is a subfield of linguistics, computer science, information engineering, and artificial intelligence concerned with the interactions between computers and human (natural) languages, in particular how to program computers to process and analyze large amounts of natural language data.



# Synonyms of NLP

- **Computational Linguistics**
- **Natural Language Processing**
- **Natural Language Understanding**
- **Human Language Processing**

- Subtleties

**Computational Linguistics** is more regarded as a branch of Linguistics, whose main purpose is to understand the mechanism of human languages by means of computing



# Synonyms of NLP

- **Computational Linguistics**
- **Natural Language Processing**
- **Natural Language Understanding**
- **Human Language Processing**

- Subtleties

**Natural Language Processing** is a branch of computer sciences and artificial intelligent, whose main purpose is to develop technologies to enable human-computer interactions using human languages



# Synonyms of NLP

- **Computational Linguistics**
- **Natural Language Processing**
- **Natural Language Understanding**
- **Human Language Processing**

- Subtleties

**Natural Language Understanding** is one of the two main challenges in Natural Language Processing, while the other is **Natural Language Generation**.



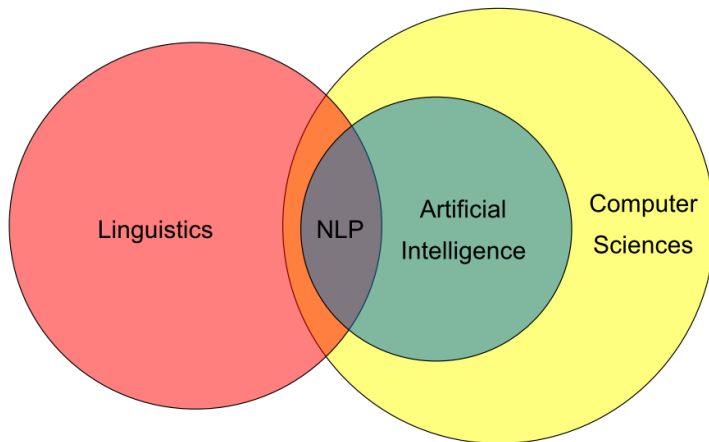
# Synonyms of NLP

- **Computational Linguistics**
- **Natural Language Processing**
- **Natural Language Understanding**
- **Human Language Processing**

- Subtleties

**Human Language Technologies** mainly refer to NLP technologies, but may also include other language related technologies, include speech technologies, optical character recognition (OCR), computer typesetting, etc.

# NLP as an interdisciplinary study



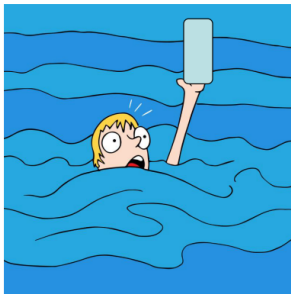
# Understanding human languages is not easy

- We are getting used to the fact that human beings can understand each other using language communication.
- Although it is a natural result of evolution for human to obtain the language competence.
- It seems to be a miracle, due to its complexity.
- No other species in this planet can use languages at the degree as humans do.
- The mechanism behind human languages is not fully discovered.
- Understanding human languages by computer is difficult.



# Understanding human languages is not easy

**Tell my wife I love her!**



**From Husband:  
I love her!**





# Research questions

- How humans understand each other by using language communication?
- Is it possible to simulate human language behaviors without understanding language mechanisms?

# The way of NLP research

- Unlike linguists who develop numerous theories to explain the language mechanisms, NLP researchers try to simulate human language behaviors by computing, not necessary to understand the language mechanisms.



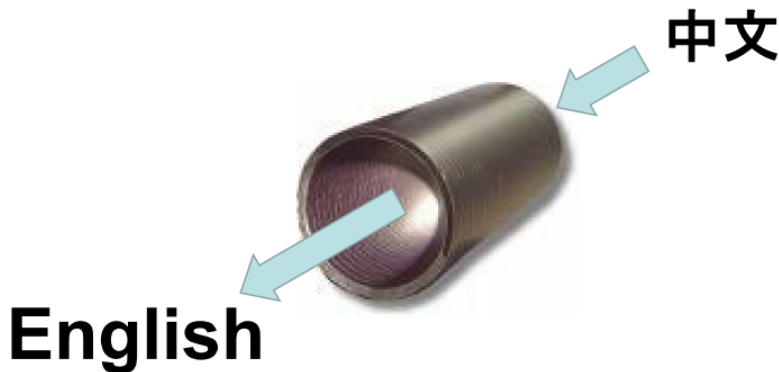
# A brief history of NLP

- 1960s-1990s: Rule-based approaches
- 1990s-2010s: Statistical approaches
- 2010s-present: Neural network (deep learning) approaches

# Holy grails of NLP

- Accurate machine translation between human languages
- Free conversation between humans and computers

# Accurate machine translation

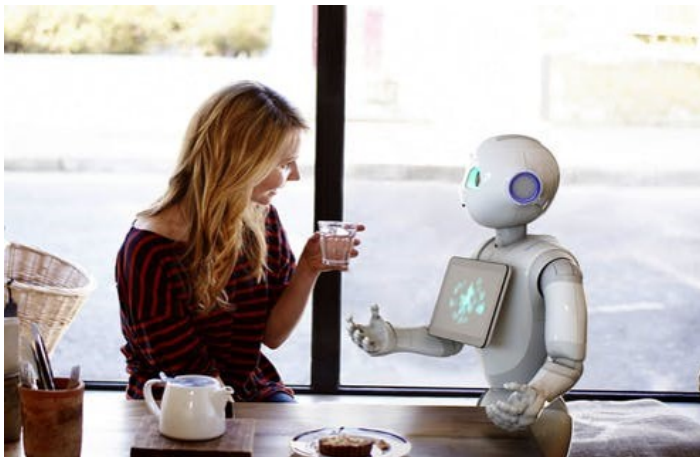


# The Tower of Babel



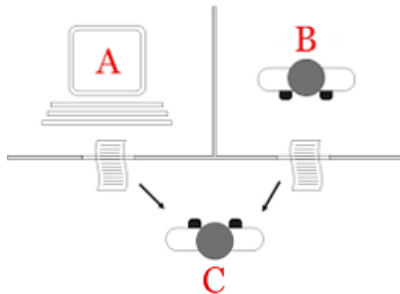
Oil painting by Pieter Bruegel the Elder, 1563, from Wikipedia

# Free human machine conversation





# Turing test

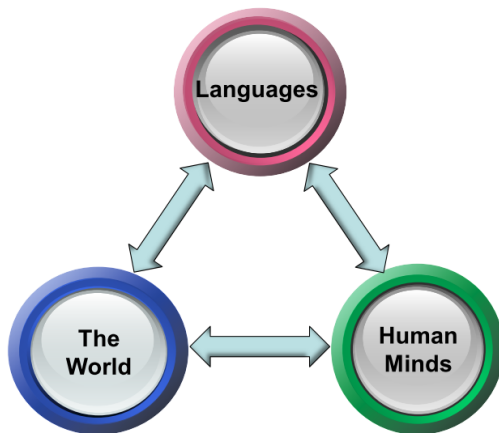


By Juan Alberto Sánchez Margallo, CC BY 2.5, from Wikipedia

# NLP Tasks

|                                  | <b>Word / Phrase</b>   | <b>Sentence</b>   | <b>Document</b>  |
|----------------------------------|--|---|--|
| Features / Expressions           | word frequency, colocations, one-hot vectors, word embeddings  | sentence embeddings, language models  | bag of word / n-grams, word frequency vectors, tf-idfs, key words / phrases extraction, topic distributions, document embeddings |
| Classification                   | part-of-speech tagging, word sense disambiguation  | sentiment analysis  | text classification, sentiment analysis  |
| Sequence labeling / Segmentation | stemming   | word segmentation, part-of-speech tagging, named entity recognition               | sentence segmentation, coreference resolution  |
| Structural prediction / Parsing  | morphological analysis   | constituent parsing, dependency parsing, semantic role labeling, semantic parsing | discourse parsing  |
| Sequence Generation /            | machine translation, text summarization, dialog, style transfer, question answering, machine reading comprehension |   |  |

# Languages, human minds and the world

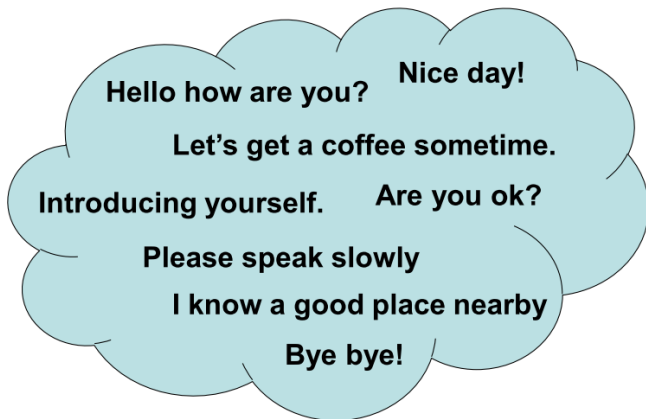




# Content

- 1 About the course
- 2 Research questions and NLP tasks
- 3 Grammars and Automata**
- 4 Text segmentation and morphology analysis
- 5 Word frequency and collocations

# How can we define a language?





# How can we define a language?

- A language can be defined as the set of sentences which can be accepted by the speakers of that language.
- It is not possible to define a natural language by enumerate all the sentences, because the number of sentences in a natural languages is infinite.
- Two feasible ways to define a language with infinite sentences:
  - By a Grammar
  - By an Automaton



# Define a language by a grammar

- A grammar  $G$  is defined as:
  - A finite set of rules, and,
  - A mechanism to generate word sequences by applying the rules in  $G$  in a finite number of time steps.
- A sentence of a language is defined as:
  - A word sequence  $S$  is called a sentence of a language  $L$  if and only if  $S$  belongs to  $L$ .
- Given a grammar  $G$ , a language  $L$  could be defined by  $G$  as:
  - A word sequence  $S$  is a sentence of  $L$  if and only if  $S$  can be generated by  $G$ .



# Define a language by an automaton (1)

- An automaton  $A$  is a abstract machine which:
  - Takes a symbol sequence  $S$  as input, and determines if  $A$  will *accept* or *reject*  $S$ .
  - Has a finite number of states and a finite number of actions.
  - At each time step,  $S$  is in a state, and points to a position in  $S$ .
  - The current state and current symbol determines the action which  $A$  will execute, which determines the next state of  $A$  and the next position of  $S$  where  $A$  will point to.
  - Given a input  $S$ ,  $A$  will run until it stops, and the final state of  $A$  determines if  $A$  will *accept* or *reject*  $S$ .

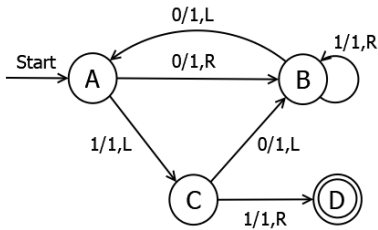
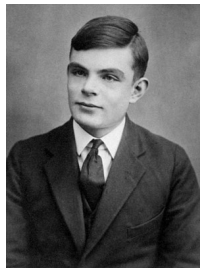
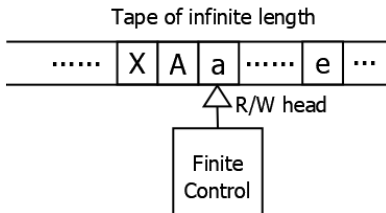




# Define a language by an automaton (2)

- A language  $L$  can be defined by an automaton  $A$  as:
  - A word sequence  $S$  is a sentence of  $L$ , if and only if: when we input  $S$  to  $A$ ,  $A$  will stop in a finite number of time steps at an *accept* state.

# Turing machine





# Turing machine

A Turing machine consists of: (to be continued)

- A *tape* divided into cells, one next to the other. Each cell contains a symbol from some finite alphabet. The alphabet contains a special blank symbol and some other symbols. The *tape* is assumed to be arbitrarily extendable to the left and to the right.
- A *read/write head* that can read and write symbols on the *tape* and move the *tape* left and right one (and only one) cell at a time.



# Turing machine

- A Turing machine consists of: (continued)
  - A *state register* that stores the state of the Turing machine, one of finitely many. Among these is the special *start state* with which the state register is initialized.
  - A finite *table* of instructions that, given the *state* the machine is currently in and the symbol it is reading on the *tape*, tells the machine to do the following in sequence:
    - Either erase or write a symbol.
    - Move the *head* to the left or right cell.
    - Assume the same or a *new state* as prescribed.



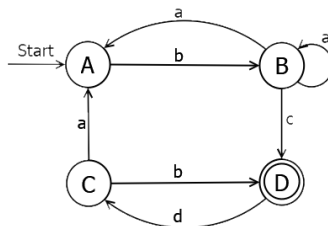
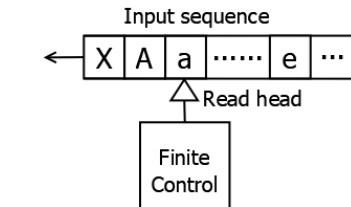
# Linear bounded automaton

A linear bounded automaton is a Turing Machine that satisfies the following three conditions:

- Its input alphabet includes two special symbols, serving as *left and right endmarkers*.
- Its *transitions* may not print other symbols over the *endmarkers*.
- Its *transitions* may neither move to the left of the *left endmarker* nor to the right of the *right endmarker*.



# Finite state automaton / machine (FSA/FSM)





# Finite state automaton / machine (FSA/FSM)

A Finite State Automaton (FSA), or Finite State Machine (FSM), consists of:

- A finite number of *states*, while the FSM can be in one *states* at each given time;
- A *head* which read a symbol from a sequence of symbols as the *input*. The *head* always goes to the next symbol at the next time step;
- A *transition* matrix which determines the next *states* of the FSM according to the current *states* and the current symbol.



# Content

- 1 About the course
- 2 Research questions and NLP tasks
- 3 Grammars and Automata
- 4 Text segmentation and morphology analysis**
- 5 Word frequency and collocations





# Text segmentation

- In NLP, text is segmented into units of various granularities, which include:
  - Chapters and sections;
  - Paragraphs;
  - Sentences;
  - Clauses;
  - Phrases;
  - Words;
  - Morphemes (stems, suffixes, prefixes).



# Text segmentation

- Text segmentation is not straightforward in many cases:
  - For languages like Chinese, Japanese, Tibetan, Thai, there are no spaces between words;
  - For languages like Thai and Tibetan, the delimiters between sentences, clauses or phrases are not ambiguous, which makes it hard to segment sentences;
  - Even for English, sentence segmentation is not a trivial task, because the full stop mark (.) is also used for abbreviations, decimals, etc., which may or may not terminate a sentence.

# Thai

โลกเราเป็นอะไรหนอในช่วงนี้ ฝั่งหนึ่งของโลกมีอากาศ  
อันแปรปรวนวิปริต หนาวเหน็บอย่างไม่เคยเกิดขึ้นมาก่อน  
และยังเกิดแผ่นดินพิโรธโกรธาคร่าชีวิตคนไปเป็นเรือนแสน  
ส่วนบ้านเรานั้นในปีที่ผ่านมาแทบไม่มีฤดูหนาวให้ชื่นใจกันเลย  
อากาศกลับร้อน แดดมีทั้งฝนหลงฤดูในช่วงนี้อีกต่างหาก  
ทุกคนพูดว่า เป็นเพราะภาวะโลกร้อนนั่นเองที่ทำให้ทุกอย่าง  
ดูไม่เหมือนเดิม ประเทศที่มีอากาศหนาวก็หนาวสุดขั้ว ประเทศ

Spaces are not reliable boundaries between sentences.

# Chinese

## 西游记 4 真假猴王

师徒四人继续西行。有一天，他们来到一个地方，前面是望不到边的水面，唐僧发愁 (chóu) 道：“这么大的水，怎么过去呢？”

四个人正不知道怎么办，忽然看见远处好像有一个人在河边，于是就走过去，想问一问。

走近了一看，那不是一个人，而是一块石头，石头上写着三个大字“通天河”，旁边还有一行小字——“河宽 (kuān) 八百里，自古少人行”，意思是这条河有八百里宽，很少有人能通过。

There are not spaces between words.



# English sentence segmentation

- Dot marks (.) are ambiguous:
  - Full stop: *This is an apple.*
  - Decimal: *235.6*
  - Abbreviations: *U.S. Ph.D. etc.*
  - A dot mark can take multiple roles: *He comes from U.S.*
- To segment English text into sentences, we need to determine whether a dot mark is an end of sentence or not.
- It can be solved as a classification problem.

# English sentence segmentation

— as a classification task

He comes from U.S. She comes from Australia.

↑   ↑   ↑

No Yes Yes

He comes from U.S. with his friends.

↑   ↑   ↑

No No Yes

# Chinese word segmentation

|     |   |  |
|-----|---|--|
| (a) | <p>下雨天留客天留我不留</p> <p>下雨、天留客。天留、我不留！</p> <p>下雨天、留客天。留我不？留！</p> | <p>Unpunctuated Chinese sentence</p> <p><i>It is raining, the god would like the guest to stay. Although the god wants you to stay, I do not!</i></p> <p><i>The rainy day, the staying day. Would you like me to stay? Sure!</i></p> |
| (b) | <p>我喜欢新西兰花</p> <p>我 喜欢 新西兰 花</p> <p>我 喜欢 新 西兰花</p>            | <p>Unsegmented Chinese sentence</p> <p><i>I like New Zealand flowers</i></p> <p><i>I like fresh broccoli</i></p>   |

<http://what-when-how.com/how-to-build-a-digital-library/word-segmentation-and-sorting-digital-library/>

Chinese word segmentation may results in different meanings.

# Chinese word segmentation

— as a character tagging task

|     |        |     |   |            |   |   |     |
|-----|--------|-----|---|------------|---|---|-----|
| S   | S      | B   | E | B          | M | E | S   |
| 我   | 有      | 一   | 台 | 计          | 算 | 机 | 。   |
| (I) | (have) | (a) |   | (computer) |   |   | (.) |

Wang & Xu, Convolutional Neural Network with Word Embeddings for Chinese Word Segmentation, IJCNLP 2017

Tags:

- **S**: single character word
- **B**: beginning character of a word
- **M**: middle character of a word
- **E**: end character of a word





# English word segmentation - Tokenization

— A example of Stanford Tokenizer

## Input

Another **ex-Golden Stater**, Paul Stankowski from **Oxnard**, is contending for a berth on the **U.S.** Ryder Cup team after winning his first PGA Tour event last year and staying within three strokes of the lead through three rounds of last **month's U.S. Open**. **H.J.** Heinz Company said it completed the sale of its **Ore-Ida** frozen-food business catering to the service industry to McCain Foods Ltd. for about **\$500** million. **It's** the first group action of its kind in Britain and one of only a handful of lawsuits against tobacco companies outside the **U.S.**

**Note:** **Text in red:** change, **text in blue:** Keep



# English word segmentation - Tokenization

— A example of Stanford Tokenizer

## Output

Another **ex-Golden Stater** , Paul Stankowski from **Oxnard** , is contending for a berth on the **U.S.** Ryder Cup team after winning his first PGA Tour event last year and staying within three strokes of the lead through three rounds of last **month's U.S. Open** . **H.J.** Heinz Company said it completed the sale of its **Ore-Ida** frozen-food business catering to the service industry to McCain Foods Ltd. for about **\$ 500** million . **It's** the first group action of its kind in Britain and one of only a handful of lawsuits against tobacco companies outside the **U.S.** .

**Note:** **Text in red:** change, **text in blue:** Keep



# Morphological analysis

- To break word down into component morphemes and build a structured representation
- A morpheme is the minimal meaning-bearing unit in a language.
  - **Stem**: the morpheme that forms the central meaning unit in a word
  - **Affix**: prefix, suffix, infix, circumfix
    - **Prefix**: e.g., possible → **im**possible
    - **Suffix**: e.g., walk → walking**ing**
    - **Infix**: e.g., hingi → h**um**ingi (Tagalog)
    - **Circumfix**: e.g., sagen → **ge**sagt**t** (German)

a slide from UW LING 570 by Fei Xia

# Two slightly different tasks

- Stemming:
  - Ex: writing  $\rightarrow$  writ + ing
- Lemmatization:
  - Ex1: writing  $\rightarrow$  write +V +Prog
  - Ex2: books  $\rightarrow$  book +N +Pl
  - Ex3: writes  $\rightarrow$  write +V +3Per +Sg



# Ambiguity in morphology

- flies  $\rightarrow$  fly +N +PL
- flies  $\rightarrow$  fly +V +3rd +Sg

a slide from UW LING 570 by Fei Xia



# Language variation

- Analytic languages: e.g., Chinese; English as a language with analytic tendency.
- Synthetic flexive languages: e.g., Russian
- Synthetic agglutinate languages: e.g., Turkish

# Ways to combine morphemes to form words

- Inflection: stem + gram. morpheme → same class
  - Ex: help + ed → helped
- Derivation: Derivation: stem + gram. morpheme → different class
  - Ex: civil + -ization → civilization
- Compounding: multiple stems
  - Ex: cabdriver, doghouse
- Cliticization: stem + clitic
  - Ex: they'll, she's (\*I don't know who she is)

a slide from UW LING 570 by Fei Xia

# UniMorph 2.0: Universal Morphology

| ARABIC                                     |                                   |                                   |                                  |                                   |                                  |                                   |                                   |                                  |  |
|--|-----------------------------------|-----------------------------------|----------------------------------|-----------------------------------|----------------------------------|-----------------------------------|-----------------------------------|----------------------------------|--|
| active voice<br>الفعل المتكلم              |                                   |                                   |                                  |                                   |                                  |                                   |                                   |                                  |  |
|  | singular<br>المفرد                |                                   |                                  | dual<br>المثنى                    |                                  |                                   | plural<br>الجمع                   |                                  |  |
|  | 1 <sup>st</sup> person<br>المتكلم | 2 <sup>nd</sup> person<br>المخاطب | 3 <sup>rd</sup> person<br>الغائب | 2 <sup>nd</sup> person<br>المخاطب | 3 <sup>rd</sup> person<br>الغائب | 1 <sup>st</sup> person<br>المتكلم | 2 <sup>nd</sup> person<br>المخاطب | 3 <sup>rd</sup> person<br>الغائب |  |
| past (perfect) indicative<br>الماضي        | m<br>رَوَيْتَ<br>zawaytu          | zawayta                           | zawwa                            | zawaytuma                         | zawayta                          | zawaytuma                         | zawaytuma                         | zawaytuma                        |  |
|  | f<br>رَوَيْتَ<br>zawaytu          | zawayta                           | zawwa                            | zawaytuma                         | zawayta                          | zawaytuma                         | zawaytuma                         | zawaytuma                        |  |
| non-past (imperfect) indicative<br>المضارع | m<br>أَرَوِي<br>azawi             | azaw                              | azaw                             | azaw                              | azaw                             | azaw                              | azaw                              | azaw                             |  |
|  | f<br>أَرَوِي<br>azawi             | azaw                              | azaw                             | azaw                              | azaw                             | azaw                              | azaw                              | azaw                             |  |

| NAVAJO       |          |           |           |
|--------------|----------|-----------|-----------|
| IMPERFECTIVE | singular | duoplural | plural    |
| 1st person   | ahleeh   | ahleeh    | da'ahleeh |
| 2nd person   | ahleeh   | ahleeh    | da'ahleeh |
| 3rd person   | ahleeh   | ahleeh    | da'ahleeh |
| 4th person   | ahleeh   | ahleeh    | da'ahleeh |
| Unspecified  | -        | Passive A | Passive B |
| Spatial      | -        | ahleeh    | -         |
| PERFECTIVE   | singular | duoplural | plural    |
| 1st person   | ahleeh   | ahleeh    | ahleeh    |
| 2nd person   | ahleeh   | ahleeh    | ahleeh    |
| 3rd person   | ahleeh   | ahleeh    | ahleeh    |
| 4th person   | ahleeh   | ahleeh    | ahleeh    |
| Unspecified  | -        | Passive A | Passive B |
| Spatial      | -        | ahleeh    | -         |

| Language | Lemma | Inflection | Features        |
|----------|-------|------------|-----------------|
| Navajo   | ateeh | da'ateeh   | V;3;PL;IPFV     |
| Arabic   | زوى   | يَرْوُونَ  | V;3;PL;IPFV;ACT |

Kirov et al., UniMorph 2.0: Universal Morphology, LREC 2018

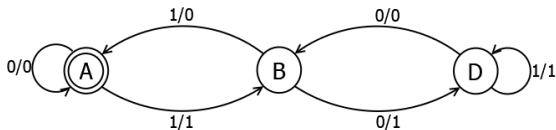




# Finite state transducers (FSTs)

- Finite State Transducers are an extension to Finite State Machines, where an output symbol will be given for each input symbol.
- FSTs are commonly used tools for morphological analysis.
- A FST can be used in a inverse direction with the input and the output swapped.

# Finite state transducers (FSTs)



| input | output |
|-------|--------|
| 0     | 0      |
| 11    | 01     |
| 110   | 010    |
| 1001  | 0011   |
| 1100  | 0100   |
| 1111  | 0101   |
| 10010 | 00110  |



# English morphology

- Affixes: prefixes, suffixes; no infixes, no circumfixes.
- Inflectional:
  - Noun: -s
  - Verbs: -s, -ing, -ed, -ed
  - Adjectives: -er, -est
- Derivational:
  - Ex:  $V + \text{suf} \rightarrow N$   
computerize + -ation  $\rightarrow$  computerization  
kill + er  $\rightarrow$  killer
- Compound: pickup, database, heartbroken, etc.
- Cliticization: 'm, 've, 're, etc.

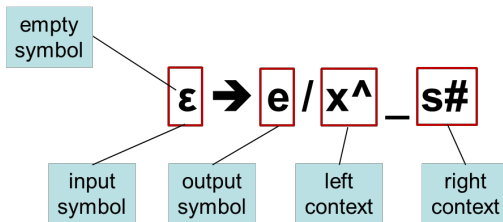
a slide from UW LING 570 by Fei Xia

# Three components

- Lexicon: the list of stems and affixes, with associated features.
  - Ex1: book: N
  - Ex2: -s: +PL
- Morphotactics:
  - Ex: +PL follows a noun
- Orthographic rules (spelling rules): to handle exceptions that can be dealt with by rules.
  - Ex3:  $\epsilon \rightarrow e / x^{\wedge} \_ s\#$

a slide from UW LING 570 by Fei Xia

# Rewrite rules



## An example

Task: foxes → fox +N +PL

Surface: foxes



Orthographic rules

Intermediate: fox ^s

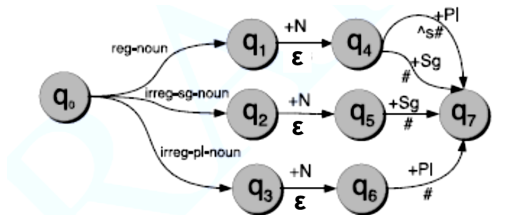


Lexicon + morphotactics

Lexical: fox +N +pl

a slide from UW LING 570 by Fei Xia

# An FST

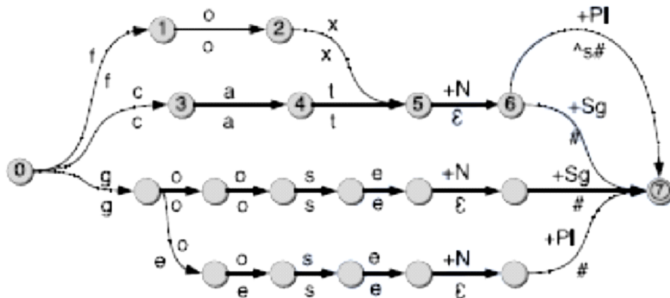


cat +N +PL  $\rightarrow$  cat  $\wedge$ s #

cat +N +Sg  $\rightarrow$  cat #

a slide from UW LING 570 by Fei Xia

# Expanding FST



fox +N + Pl  $\rightarrow$  fox  $\wedge$  s #

cat +N + Pl  $\rightarrow$  cat  $\wedge$  s #

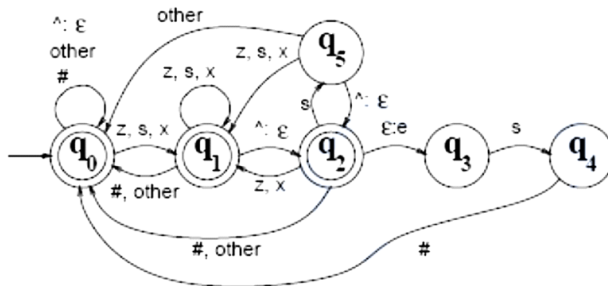
goose +N +Sg  $\rightarrow$  goose #

goose +N +Pl  $\rightarrow$  geese #

a slide from UW LING 570 by Fei Xia



# Representing orthographic rules as FSTs



$\epsilon \rightarrow e / (s|x|z) \hat{\_} s \#$

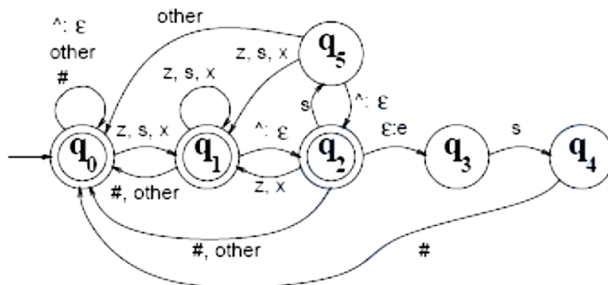
Input:  $\dots(s|x|z) \hat{s} \#$  immediate level

Output:  $\dots(s|x|z)es \#$  surface level

To reject (fox  $\hat{s}$ , foxs)

a slide from UW LING 570 by Fei Xia

# Representing orthographic rules as FSTs



(fox, fox):  $q_0, q_0, q_0, q_1$

(fox#, fox#):  $q_0, q_0, q_0, q_1, q_0$

(fox<sup>z</sup>#, foxz#):  $q_0, q_0, q_0, q_1, q_2, q_1, q_0$

(fox<sup>s</sup>#, foxes#):  $q_0, q_0, q_0, q_1, q_2, q_3, q_4, q_0$

(fox<sup>s</sup>, foxs):  $q_0, q_0, q_0, q_1, q_2, q_5$

a slide from UW LING 570 by Fei Xia

# Further reading on morphological analysis

- Fei Xia, slides on morphological analysis

[https://www.powershow.com/viewfl/6a39a-ZDc1Z/Morphological\\_analysis\\_powerpoint\\_ppt\\_presentation](https://www.powershow.com/viewfl/6a39a-ZDc1Z/Morphological_analysis_powerpoint_ppt_presentation)

- Mans Hulden (2011), Morphological analysis with FSTs

<https://fomafst.github.io/morptut.html>



# Content

- 1 About the course
- 2 Research questions and NLP tasks
- 3 Grammars and Automata
- 4 Text segmentation and morphology analysis
- 5 Word frequency and collocations**



# Top 5000 words in American English

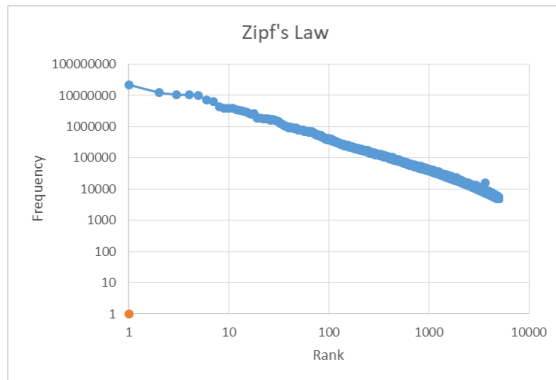
| Rank | Word | Part of speech | Frequency | Dispersion |
|------|------|----------------|-----------|------------|
| 1    | the  | a              | 22038615  | 0.98       |
| 2    | be   | v              | 12545825  | 0.97       |
| 3    | and  | c              | 10741073  | 0.99       |
| 4    | of   | i              | 10343885  | 0.97       |
| 5    | a    | a              | 10144200  | 0.98       |
| 6    | in   | i              | 6996437   | 0.98       |
| 7    | to   | t              | 6332195   | 0.98       |
| 8    | have | v              | 4303955   | 0.97       |
| 9    | to   | i              | 3856916   | 0.99       |
| 10   | it   | p              | 3872477   | 0.96       |

| Rank | Word | Part of speech | Frequency | Dispersion |
|------|------|----------------|-----------|------------|
| 1    | the  | a              | 22038615  | 0.98       |
| 2    | be   | v              | 12545825  | 0.97       |
| 3    | and  | c              | 10741073  | 0.99       |
| 4    | of   | i              | 10343885  | 0.97       |
| 5    | a    | a              | 10144200  | 0.98       |
| 6    | in   | i              | 6996437   | 0.98       |
| 7    | to   | t              | 6332195   | 0.98       |
| 8    | have | v              | 4303955   | 0.97       |
| 9    | to   | i              | 3856916   | 0.99       |
| 10   | it   | p              | 3872477   | 0.96       |

Statics from Corpus of the Contemporary American English

<http://www.wordfrequency.info/>

# Top 5000 words in American English



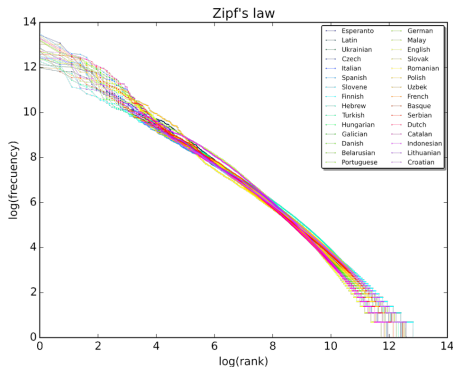


# Zipf's Law

The frequency of any word is inversely proportional to its rank in the frequency table:

$$p(w_r) \propto \frac{1}{r}$$

# Zipf's law



A plot of the rank versus frequency for the first 10 million words in 30 Wikipedias (dumps from October 2015) in a log-log scale.

(By SergioJimenez - Own work, CC BY-SA 4.0, from Wikipedia)





# Collocation or multi-word expression (MWE)

- A COLLOCATION is an expression consisting of two or more words that correspond to some conventional way of saying things.
- The words together can mean more than their sum of parts
  - The Times of India, disk drive
  - hot dog, mother in law



# Collocation or multi-word expression (MWE)

- Examples of collocations
  - noun phrases like *strong tea* and *weapons of mass destruction*
  - phrasal verbs like to *make up*, and other phrases like the *rich and powerful*.
- Valid or invalid?
  - *a stiff breeze* but not a *stiff wind* (while either a *strong breeze* or a *strong wind* is okay).
  - *broad daylight* (but not bright daylight or narrow darkness).

Manning & Schütze, Fundamentals of Statistical Natural Language Processing, 1999



# Criteria for collocations (or MWE)

- Typical criteria for collocations:
  - non-compositionality
  - non-substitutability
  - non-modifiability.
- Collocations usually cannot be translated into other languages word by word.
- A phrase can be a collocation even if it is not consecutive (as in the example *knock ... door*).

Manning & Schütze, Fundamentals of Statistical Natural Language Processing, 1999



# Non-Compositionality

- A phrase is compositional if the meaning can be predicted from the meaning of the parts.
  - E.g. new companies
- A phrase is non-compositional if the meaning cannot be predicted from the meaning of the parts
  - E.g. *hot dog*

Manning & Schütze, Fundamentals of Statistical Natural Language Processing, 1999



# Non-Compositionality

- Collocations are not necessarily fully compositional in that there is usually an element of meaning added to the combination.
  - E.g. *strong tea*
- Idioms are the most extreme examples of non-compositionality
  - E.g. *to hear it through the grapevine*

Manning & Schütze, Fundamentals of Statistical Natural Language Processing, 1999



# Non-Substitutability

- We cannot substitute near-synonyms for the components of a collocation.
- For example
  - We can't say *yellow wine* instead of *white wine* even though *yellow* is as good a description of the color of *white* wine as white is (it is kind of a yellowish white).

Manning & Schütze, Fundamentals of Statistical Natural Language Processing, 1999



# Non-Substitutability

- Many collocations cannot be freely modified with additional lexical material or through grammatical transformations (Non-modifiability).
  - E.g. *white wine*, but not *whiter wine*
  - E.g. *mother in law*, but not *mother in laws*

Manning & Schütze, Fundamentals of Statistical Natural Language Processing, 1999



# Metrics for Collocation or MWE Extraction

- Frequency
- Mean and Variance of Distances between Words
- Hypothesis Testing
  - $t$ -test
  - $\chi^2$  test
  - likelihood ratio test
- Mutual Information
- Left and Right Context Entropy
- C-Value





# Further reading on collocation and MWE

- Manning & Schütze, Fundamentals of Statistical Natural Language Processing, 1999, Chapter 3 (A general introduction to collocation)
- Katerina T. Frantzi, Sophia Ananiadou, Junichi Tsujii, The C-value / NC-value Method of Automatic Recognition for Multi-word Terms, ECDL 1998: Research and Advanced Technology for Digital Libraries pp 585-604 (proposed the C-value metric)
- Zhiyong Luo, Rou Song, An integrated method for Chinese unknown word extraction, SIGHAN 2004. Barcelona, Spain. (proposed the context entropy method)



# Content

- 1 About the course
- 2 Research questions and NLP tasks
- 3 Grammars and Automata
- 4 Text segmentation and morphology analysis
- 5 Word frequency and collocations