

## Basic Notions and Low-Level Text Normalization

## 1. Basic Notions: Characters, Words, Tokens, Documents, Corpora

- a. The human race has developed many different writing systems, based on several categories of graphemes (atomic symbols). To vastly over-simplify, we have
  - i. Alphabets: a set of < 100 symbols, each roughly corresponding to a speech sound:

## 1. Example

- a. English: a b c ... z
- b. Greek: alpha, beta, ...

- ii. Syllabaries: a set of 100s of symbols, each roughly corresponding to spoken syllable

## 1. Example

Linear B (early Greek):    †    ॥    pa-te = πατήρ = pater = "father"

- a. Japanese: さか, saka, "hill"

- iii. Logographies: A set of 1000s of symbols, each roughly corresponding to a spoken word or concept

## 1. Example

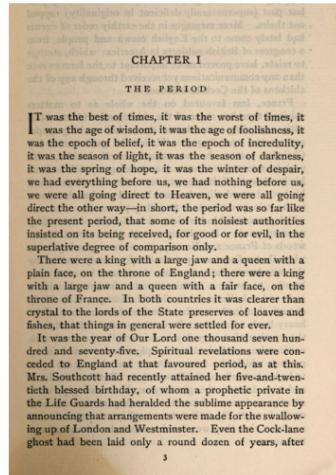
Egyptian Hieroglyphs:                

Chinese: 山 · mountain. 刀 knife 下 down

## a.

- b. For NLP, we want to process text as a sequence of atomic symbols and these may be any of the preceding categories
- c. Thus: In NLP, textual data is presented in its most basic form as a sequence of atomic symbols from some finite collection (think Unicode!).
  - i. Unicode is more modern version
- d. In CS 505, our language is English, and this collection will be ASCII characters, and we will generally just call them characters. Thus, in its most basic form, a text is simply one long string.
- e. This string form of a text is the minimal representation of the information content of the text (excluding formatting, diagrams, different fonts, illustrations, etc.) and may include some minimal formatting (white space, \n, \t, etc.):

## A Tale of Two Cities, Charles Dickens



f.

- g. Although we will have occasion to use the string form when we study character-level machine learning models, almost all NLP uses data which has been grouped into larger units:
  - i. Words: Sequence of characters, separated by white space or punctuation
  - ii. Tokens: Words possibly preprocessed into some more useful form (the rest of this lecture)
  - iii. Sentences: Sequences of words/tokens
  - iv. Paragraphs: Sequences of sentences
  - v. Chapters/Sections/Topics: Sequences of paragraphs
  - vi. Document: Sequence of paragraphs
  - vii. Corpus: Set of documents

## 2. Corpora for NLP

- a. There are many publicly-available corpora for NLP, often categorized (and preprocessed) for specific tasks, in various languages, etc. Dr Google will help you find these....

**1.1 Gutenberg Corpus**  
NLTK includes a small selection of texts from the Project Gutenberg electronic text archive, which contains over 20,000 books. To use this corpus, run `nltk.corpus.gutenberg.readme()`.

**1.3 Brown Corpus**  
The Brown Corpus is a collection of about 500 million words of text from 500 different sources, mostly from American English.

**1.4 Reuters Corpus**  
The RCV1 corpus is a collection of news stories from the Reuters newswire, with approximately 800,000 documents.

**1.5 Inaugural Address Corpus**  
In 1789, George Washington gave his first inaugural address. This corpus contains the text of every presidential inauguration since that time.

**1.7 Corpora in Other Languages**  
NLTK comes with corpora for many languages, though in some cases you will need to download them separately. These include:

```
>>> nltk.corpus.cess_esp.words()
['El', 'grupo', 'estatal', 'Electricit\u00e9s de France', ...]
['Un', 'revolucionario', 'referentes', 'O', 'de', 'M\u00e9jico', ...]
>>> nltk.corpus.esperanto.words()
['V', 'algun', 'se', 'que', 'quier', '...']
```

**Common Crawl** maintains a **free repository of web data** that can be used for AI training datasets.

**Enron Corpus**  
The Enron Corpus is a database of over 200,000 emails from the years leading up to the company's collapse in December 2001. It is a collection of over 200,000 emails from the Enron email system, which contained over 500,000 individual messages. The emails are in plain text format and are available for download.

Dataset	Raw Size	Weight	Composition	Amplification
CommonCrawl	410 B	60%	82%	-0.27
RCV1v2	19 B	20%	4%	+5.00
Books1.8	67 B	15%	13%	+0.15
Books2	3 billion	5%	1%	+5.00
<b>TOTAL</b>	<b>500B tokens</b>			

b.

**Datasets**

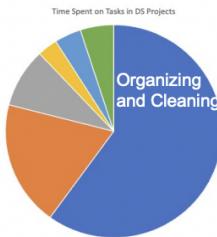
**bookcorpus** 137

Tasks: Text Generation, Fill-Mask, Sub-tasks: language-modeling, masked-language-modeling

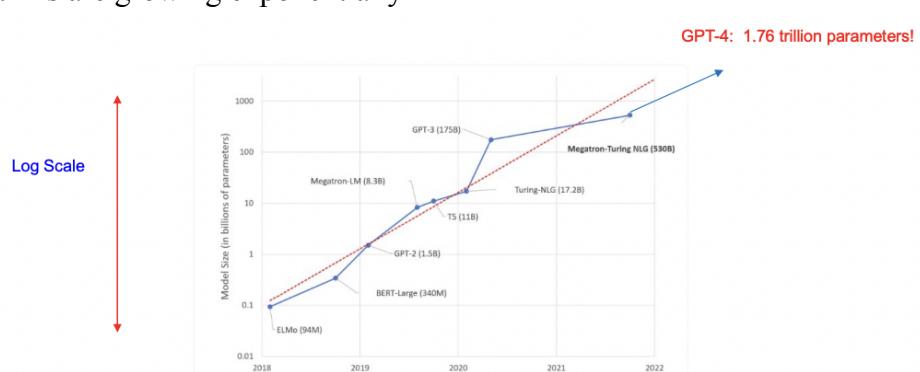
Annotations Creators: no-annotation Source Datasets: original Author: arXiv License: CC-BY-NC-ND

### 3. Textual Data Preparation

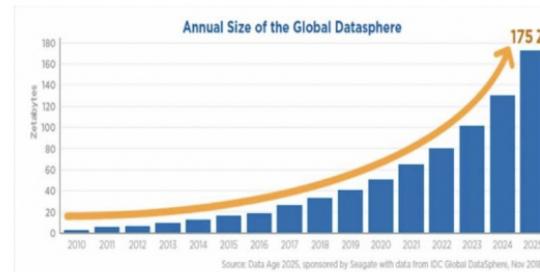
- a. Why do we need to learn low-level text processing?
  - i. Because these corpora are for education, contests, creating general language models (e.g., chatGPT), etc. Most NLP projects involve taking some raw textual data and wrangling it into a corpus of your own.
- b. CrowdFlower, provider of a “data enrichment” platform for data scientists, conducted a survey of about 80 data scientists and found that data scientists spend
  - i. 60% of the time in organizing and cleaning data
  - ii. 19% of the time is spent in collecting datasets
  - iii. 9% of the time is spent in mining the data to draw patterns
  - iv. 3% of the time is spent in training the datasets
  - v. 4% of the time is spent in refining the algorithms
  - vi. 5% of the time is spent in other tasks



- c. ML algorithms are growing exponentially

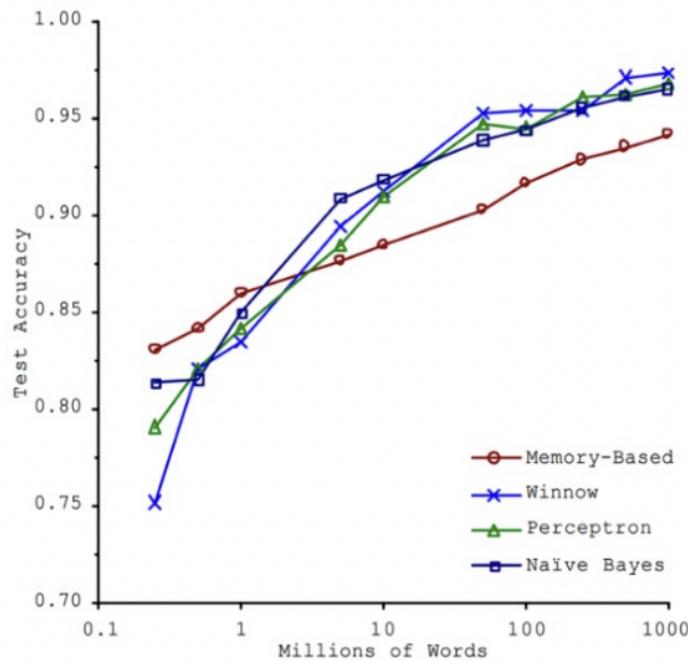


- i.
- d. Basic rule: the more parameters, the more data you need
- e. The amount of available data is growing exponentially, and it is mostly unstructured. It is simply impossible to process this tsunami of data by (human) hand!



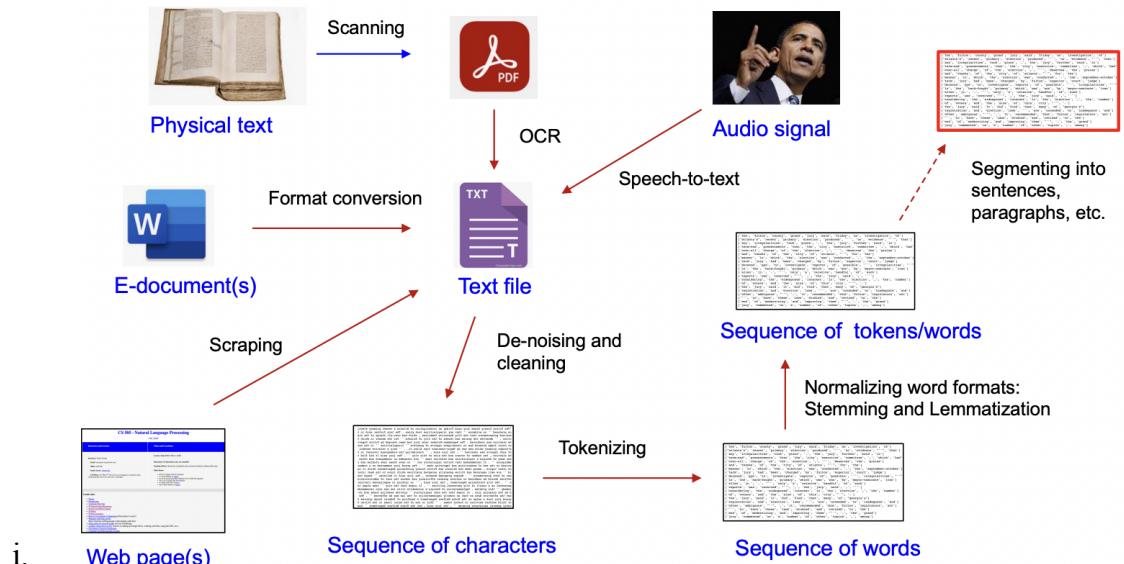
- i.

- f. Plus, data is the most important resource – progress in NLP is overwhelmingly dependent on the amount of data, NOT refinements to the algorithms. (need an enormous amount of data → the more data you have, the better the solution)



i.

- g. Most textual data preparation will follow some part of this flowchart



i.

#### 4. Tokenizing: Separating a string into words

##### a. First try: separate on white space

Everybody was drunk. The whole battery was drunk going along the road in the dark. We were going to the Champagne. The lieutenant kept riding his horse out into the fields and saying to him, "I'm drunk, I tell you, mon vieux. Oh, I am so souused." We went along the road all night in the dark and the adjutant kept riding up alongside my kitchen and saying, "You must put it out. It is dangerous. It will be observed." We were fifty kilometers from the front but the adjutant worried about the fire in my kitchen. It was funny going along that road. That was when I was a kitchen corporal.

```
1 import re
2
3 separator_1 = '\s+'          # one or more white space characters
4
5 print(re.split(separator_1, text), '\n')
6
['in', 'our', 'time', 'chapter', 'l', 'Everybody', 'was', 'drunk.', 'The', 'whole', 'battery', 'was', 'drunk', 'goi
g', 'along', 'the', 'road', 'in', 'the', 'dark.', 'We', 'were', 'going', 'to', 'the', 'Champagne.', 'The', 'lieutenan
t', 'kept', 'riding', 'his', 'horse', 'out', 'into', 'the', 'fields', 'and', 'saying', 'to', 'him', "'I'm', 'drun
k', 'I', 'tell', 'you', 'mon', 'vieux.', 'Oh,', 'I', 'am', 'so', 'souused.', 'We', 'went', 'along', 'the', 'road',
'all', 'night', 'in', 'the', 'dark', 'and', 'the', 'adjutant', 'kept', 'riding', 'up', 'alongside', 'my', 'kitchen',
'and', 'saying', '"You', 'must', 'put', 'it', 'out.', 'It', 'is', 'dangerous.', 'It', 'will', 'be', 'observed.', 'W
e', 'were', 'fifty', 'kilometers', 'from', 'the', 'front', 'but', 'the', 'adjutant', 'worried', 'about', 'the', 'fir
e', 'in', 'my', 'kitchen.', 'It', 'was', 'funny', 'going', 'along', 'that', 'road.', 'That', 'was', 'when', 'I', 'wa
s', 'a', 'kitchen', 'corporal.']

i. Punctuations were included → period, comma, quotations, and etc.
```

##### b. Second try: separate on white space and punctuation

```
1 separator_2 = '[\s,;.:!?]+'      # one or more white space or punctuation characters
2
3 print(re.split(separator_2, text), '\n')
4
['in', 'our', 'time', 'chapter', 'l', 'Everybody', 'was', 'drunk', 'The', 'whole', 'battery', 'was', 'drunk', 'goi
g', 'along', 'the', 'road', 'in', 'the', 'dark', 'We', 'were', 'going', 'to', 'the', 'Champagne', 'The', 'lieutenan
t', 'kept', 'riding', 'his', 'horse', 'out', 'into', 'the', 'fields', 'and', 'saying', 'to', 'him', "'I'm', 'drun
k', 'I', 'tell', 'you', 'mon', 'vieux', 'Oh', 'I', 'am', 'so', 'souused', 'We', 'went', 'along', 'the', 'road', 'al
l', 'night', 'in', 'the', 'dark', 'and', 'the', 'adjutant', 'kept', 'riding', 'up', 'alongside', 'my', 'kitchen', 'an
d', 'saying', '"You', 'must', 'put', 'it', 'out', 'It', 'is', 'dangerous', 'It', 'will', 'be', 'observed', 'We', 'we
re', 'fifty', 'kilometers', 'from', 'the', 'front', 'but', 'the', 'adjutant', 'worried', 'about', 'the', 'fire',
'in', 'my', 'kitchen', 'It', 'was', 'funny', 'going', 'along', 'that', 'road', 'That', 'was', 'when', 'I', 'was',
'a', 'kitchen', 'corporal.']

i. Still included the quotation marks (better since there are no
periods/commas)
```

##### c. Third try: separate on anything other than a word character

```
1 separator_3 = '\W+'          # one or more non-word characters
2
3 print(re.split(separator_3, text), '\n')
4
['in', 'our', 'time', 'chapter', 'l', 'Everybody', 'was', 'drunk', 'The', 'whole', 'battery', 'was', 'drunk', 'goi
g', 'along', 'the', 'road', 'in', 'the', 'dark', 'We', 'were', 'going', 'to', 'the', 'Champagne', 'The', 'lieutenan
t', 'kept', 'riding', 'his', 'horse', 'out', 'into', 'the', 'fields', 'and', 'saying', 'to', 'him', 'I', 'm', 'drun
k', 'I', 'tell', 'you', 'mon', 'vieux', 'Oh', 'I', 'am', 'so', 'souused', 'We', 'went', 'along', 'the', 'road', 'all',
'night', 'in', 'the', 'dark', 'and', 'the', 'adjutant', 'kept', 'riding', 'up', 'alongside', 'my', 'kitchen', 'and',
'saying', '"You', 'must', 'put', 'it', 'out', 'It', 'is', 'dangerous', 'It', 'will', 'be', 'observed', 'We', 'we
re', 'fifty', 'kilometers', 'from', 'the', 'front', 'but', 'the', 'adjutant', 'worried', 'about', 'the', 'fire', 'in', 'm
y', 'kitchen', 'It', 'was', 'funny', 'going', 'along', 'that', 'road', 'That', 'was', 'when', 'I', 'was', 'a', 'kitch
en', 'corporal.']

i. Pretty clean except there are some errors → I'm becomes I & m, and etc.
```

- d. Can't just blindly remove punctuation or non-word characters
  - i. m.p.h., Ph.D., AT&T, cap'n
  - ii. prices (\$45.55)
  - iii. dates (01/02/06)
  - iv. URLs (<http://www.stanford.edu>)
  - v. hashtags (#nlproc)
  - vi. email addresses (someone@cs.colorado.edu) o
- e. Clitic: a word that doesn't stand on its own
  - i. "are" in we're, French "je" in j'ai, "le" in l'honneur o
- f. When should multiword expressions (MWE) be words?
- i. New York, rock 'n' roll
- g. Even worse, many languages do not use consistently use spaces or punctuation to separate words, and the tokenization is quite complicated!

### Chinese

莎拉波娃现在居住在美国东南部的佛罗里达。

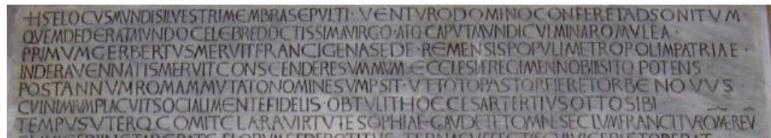
莎拉波娃 现在 居住 在 美国 东南部 的 佛罗里达

Sharapova now lives in US southeastern Florida

### Classical Sanskrit

पश्यैतां पाण्डुपुत्राणामाचार्य महतीं चमूम् ।

### Latin Inscriptions:



- h. In the case of English and many western languages, there are a variety of useful approaches
  - i. Rule-based methods, perhaps with list of exceptions
  - ii. Sequence-to-Sequence methods
    - 1. Hidden Markov Models
    - 2. Neural Networks
- i. Example: Tokenization in the NLTK uses rules based on regular expressions

```
>>> text = 'That U.S.A. poster-print costs $12.40...'
>>> pattern = r'''(?x)      # set flag to allow verbose regexps
...     ([A-Z]\.)+        # abbreviations, e.g. U.S.A.
...     | \w+(-\w+)*       # words with optional internal hyphens
...     | \$?\d+(\.\d+)?%? # currency and percentages, e.g. $12.40, 82%
...     | \.\.\.            # ellipsis
...     | [] [.,;'?():-_'] # these are separate tokens; includes ], [
...     ''
>>> nltk.regexp_tokenize(text, pattern)
['That', 'U.S.A.', 'poster-print', 'costs', '$12.40', '...']
```

- j. Example: SpaCy uses a multi-phase approach based on rules and exceptions
- First, the tokenizer split the text on whitespace similar to the split() function
  - Then, the tokenizer checks whether the substring matches the tokenizer exception rules
  - For example, “don’t” does not contain whitespace, but should be split into two tokens, “do” and “n’t”, while “U.K.” should always remain one token
  - Next, it checks for a prefix, suffix, or infix in a substring, these include commas, periods, hyphens, or quotes. If it matches, the substring is split into two tokens

```
import spacy

nlp = spacy.load("en_core_web_sm")
doc = nlp(text)
tst = [token.text for token in doc]

for k in range(0, len(tst), 10):
    print(tst[k:k+10])
```

['in', 'our', 'time', '\n\n\n', 'chapter', 'l', '\n\n', 'Everybody', 'was', 'drunk']
[.', 'The', 'whole', 'battery', 'was', 'drunk', 'going', 'along', 'the', 'road']
['in', 'the', 'dark', '.', 'We', 'were', 'going', 'to', 'the', 'Champagne']
[.', 'The', 'lieutenant', 'kept', 'riding', 'his', 'horse', 'out', 'into', 'the']
['fields', 'and', 'saying', 'to', 'him', 'I', 'm', 'drunk']
[., 'I', 'tell', 'you', 'mon', 'vieux', 'Oh', '']
['I', 'am', 'so', 'soused', 'We', 'went', 'along', 'the']
['road', 'all', 'night', 'in', 'the', 'dark', 'and', 'the', 'adjudant', 'kept']
['riding', 'up', 'alongside', 'my', 'kitchen', 'and', 'saying', 'You']
['must', 'put', 'it', 'out', 'It', 'is', 'dangerous', 'It']
['will', 'be', 'observed', 'We', 'were', 'fifty', 'kilometers', 'from']
['the', 'front', 'but', 'the', 'adjudant', 'worried', 'about', 'the', 'fire', 'in']
['my', 'kitchen', 'It', 'was', 'funny', 'going', 'along', 'that', 'road']
[., 'That', 'was', 'when', 'I', 'was', 'a', 'kitchen', 'corporal', '.']

V.

+ Code + Text

## 5. Word Normalization; Stemming and Lemmatization

- Normalization is putting words into a standard format
- Simple: Make text case-insensitive by converting all to lower case
- More complex:
  - Misspellings
    - In tge beginning....“
  - Abbreviations:
    - PHD PhD Ph.D Ph.D.
    - etc. &c
    - US U.S. U.S.A.
  - Hyphenation:
    - lowercase lower-case
  - Miscellaneous:
    - uhhuh uh-huh

- d. Stemming: remove suffixes
  - i. Not always good
  - ii. Useful when we do classification
  - iii. Not useful with translation
- e. likes, liked, likely, liking, likable     Stem: like
- f. Naïve method: Chop off last part of word (based on list of cases)!
  - i. Pretty well: chop off last 20% of the word (first part of the word was enough to distinguish)

This was not the map we found in Billy Bones's chest, but an accurate copy, complete in all things-names and heights and soundings-with the single exception of the red crosses and the written notes.

Thi wa not the map we found in Billi Bone s chest but an accur copi complet in all thing name and height and sound with the singl except of the red cross and the written note .

- ii.
- g. Better: Multi-phase rule-based systems such as the Porter Stemmer (available in NLTK)

ATIONAL → ATE (e.g., relational → relate)

ING → ε if stem contains vowel (e.g., motoring → motor)

SSES → SS (e.g., grasses → grass)

i.

ii. Does not always work but is good enough

- h. Lemmatization: Represent all words by their lemma, the shared root word (dictionary headword)

▪ *am, are, is* → *be*

▪ *car, cars, car's, cars'* → *car*

▪ Spanish *quiero* ('I want'), *quieres* ('you want')

→ *querer* 'want'

▪ *He is reading detective stories*

i. → *He be read detective story*

ii. Rules that amount to bunch of cases (better on classification)

## 6. Sentence Segmentation

- a. Segmenting into sentences is based on punctuation
- b. !, ? Are mostly unambiguous but period “.” is very ambiguous
  - i. Usually “.” is a sentence boundary. Especially followed by capital letter.
  - ii. But Abbreviations like Inc. or Dr. Snyder && Numbers like .02% or 4.3
- c. Common algorithm: Tokenize first: use rules or ML to classify a period as either (a) part of the word or (b) a sentence-boundary. Capitalization can help too!
- d. Sentence segmentation can then often be done by rules based on this tokenization.