



THE UNIVERSITY
of ADELAIDE

CRICOS PROVIDER 00123M

ANLP week1: Introduction part 1

Dr Alfred Krzywicki
University of Adelaide

adelaide.edu.au

seek LIGHT

Welcome to Week 1 of Applied NLP!

Natural Language Processing

“Natural language processing is a range of computational techniques for analyzing and representing naturally occurring texts at one or more levels of linguistic analysis for the purpose of achieving human-like languages processing for a range of particular tasks or applications.”
by Liddy (1998)

Some other names:

Computational Linguistics
Natural Language Engineering
Speech and Text Processing

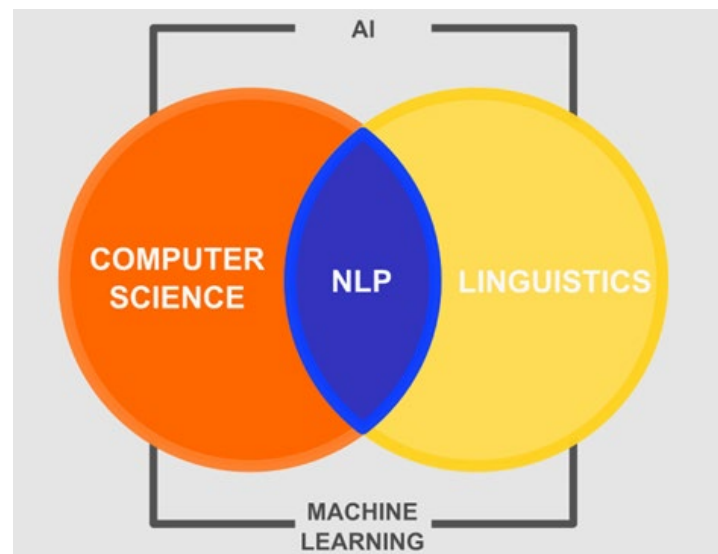


Image source: algorithmxlab.com

Language

A **vocabulary** consists of a set of **words** (w_i)



A **text** is composed of a sequence of **words** from a **vocabulary**



Beyond the genome

Studies of the epigenomic signatures of many healthy and diseased human tissues could provide crucial information to link genetic variation and disease.

The Greek prefix *epi-* can signify upon, on, over, near, at, before, and after. Most of these could apply to its use in the term 'epigenetics'—particularly the last of these. It is some 14 years, almost to the day, that Nature published the draft sequence of the human genome. Now, in this issue, we publish results from a subsequent study on the non-genomic modifications to the genome—epigenetic modifications—that critically determine which genes are turned on or off, which are turned off and which are turned on, starting on page 313, as well as in several other Nature Publishing Group journals.

Insights into these fundamental aspects of epigenetics emerge: how the epigenome affects gene expression, how the epigenome changes during stem differentiation that is, during normal development and how it changes during disease.

The results emphasize the central role of epigenetic information in the control of gene expression and cellular differentiation. *—G. G. Skellern*

It is hard to think of any branch of human biology that has not benefited from the human genome sequence. Its legacy has perhaps been most notable in advances in our appreciation of the part that genetics and genetic variations play in the normal functioning of a human body and in disease. But despite the progress, each question that the genome helps to answer throws up further questions. Much remains to be understood about how genetic information is interpreted by the

This is where epigenetics comes in. Upon the genome, on the genome, over the genome—take your pick—epigenetics collectively describes changes in the regulation of gene expression that can be passed on to a cell's progeny but are not due to changes to the nucleotide sequence of the genome.

Soon after the human genome sequence had been completed, it became clear that an epigenome — a map of the genome-wide modifications made to DNA and the protein scaffold that supports it — would also be required. The task at hand was, in researchers' like to say, not trivial. Every cell in the body carries the same genome (with a few exceptions), but the epigenome changes with cell and tissue type.

Epigenetics is still an emerging science, but researchers are now building tools to study epigenetic changes in the genome in a systematic and genome-wide way. In 2012, Nature celebrated the publication of the results of the ENCODE project, the aim of which was to describe all the functional elements encoded in the human genome by mapping

ENCODE was a pioneer in scale of effort and development of specialized analytical software, and has already had a tremendous impact on human-genetics studies. But its clinical application is limited because most of its results come from a small number of laboratory cell lines. Clinically useful epigenetic information must instead be drawn directly from all the different cell types that make up the human body.

This type of epigenomic information has now been gathered, in the Roadmap Epigenomics Project directed by the US National Institutes of Health. This project set out to generate and publicly share epigenomic data from stem cells, from mature cells from a variety of different tissues from healthy people, and from patients with diseases such as cancer, and neurodegenerative and autoimmune disease.

starting on page 313, as well as in several other Nature Publishing Group journals.

thoughts into three fundamental aspects of epigenetics emerge: how the epigenome affects gene expression; how the epigenome changes during stem-cell differentiation (that is, during normal development

The results emphasize the central role of epigenomic information in understanding these processes. Crucially,

Tackling disease using information

on the genome alone has been like trying to

like trying to work with one hand tied behind

however, if we are to understand the underlying disease mechanism and design targeted treatments. With the new wealth of data, consistent alteration in the

epigenetic landscape could identify candidate genes and pathways for further follow-up. And time-course studies of the epigenetics of cell types relevant to a specific disease could indicate whether epigenetic

One reason that it has been difficult to relate some diseases to diet is that DNA mutations (rather than those of the immune system)

disruption in tRNA function is that many of the key changes occur in poorly understood regions of the genome, usually outside those parts that code for proteins. Epigenetic maps such as those pub-

linked today should help scientists to navigate this poorly charted landscape. By overlaying these maps, made in relevant cell types, researchers can determine, for example, whether an epigenetic

change associated with a given disease lies in a region of the genome that regulates gene activity. If it does, then this overlap provides a possible lead to be explored.

Cancer is often called the disease of the genome, but the genome does not exist, or operate, in splendid isolation. Of all diseases, can-

cer has been linked most unambiguously to epigenetic aberrations. Scientists have long suspected that epigenomic organization affects the genomic location of the mutations that provoke cancer. The new

findings suggest that this is true, and they go further. They show that the epigenome of a cancer cell carries a fingerprint of the cell type that originated the cancer. This is crucial information, especially in

In human diseases, the disease and environment interact together

Tackling disease using information on the genome alone has been like trying to work with one hand tied behind the back. The new tools

of epigenomics, data fires the other hand. It will not provide all the answers. But it could help researchers decide which questions to ask.

19 FEBRUARY 2012 | VOL. 510 | NATURE | 27
© 2012 Macmillan Publishers Limited. All rights reserved



(<http://www.old-engli.sh/language.php>)

Artificial Language

```
try {
    cMessage = messageQueue.take();
    for (AsyncContext ac : queue) {
        try {
            PrintWriter acWriter = ac.getWriter();
            acWriter.println(cMessage);
            acWriter.flush();
        } catch (IOException e) {
            System.out.println("Error: " + e.getMessage());
        }
    }
} catch (InterruptedException e) {
    e.printStackTrace();
}
```

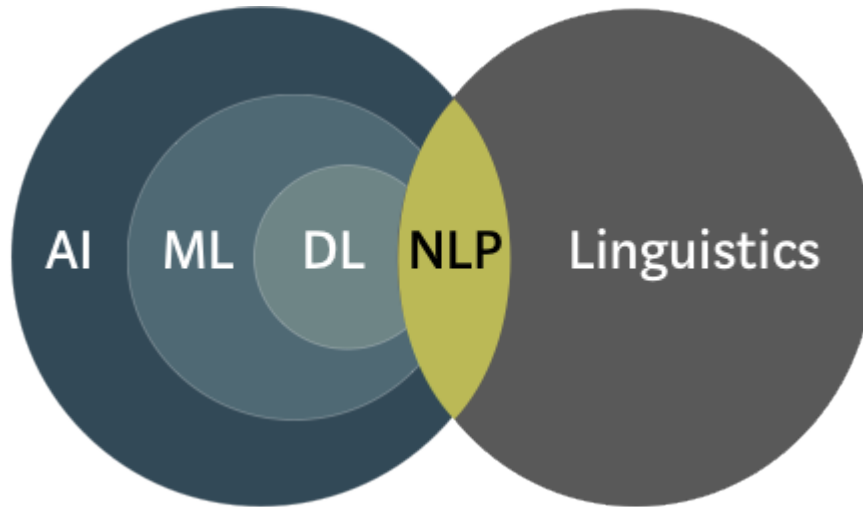
(<https://netbeans.org/features/java/>)

```
def add5(x):
    return x+5

def dotwrite(ast):
    nodename = getNodeName()
    label=symbol.sym_name.get(int(ast[0]),ast[0])
    print '    %s [label="%s" % (nodename,label),
    if isinstance(ast[1], str):
        if ast[1].strip():
            print '= %s';' % ast[1]
        else:
            print ''
    else:
        print '[';
        children = []
        for n, child in enumerate(ast[1:]):
            children.append(dotwrite(child))
        print '    %s -> (' % nodename,
        for name in children:
            print '%s' % name,
```

(<http://noobite.com/learn-programming-start-with-python/>)

Natural Language Processing



Natural language processing (NLP): how to program computers to process and analyse large amounts of natural language data.

Why NLP is hard?

- Natural language is highly ambiguous

“At last, a computer that understands you like your mother”

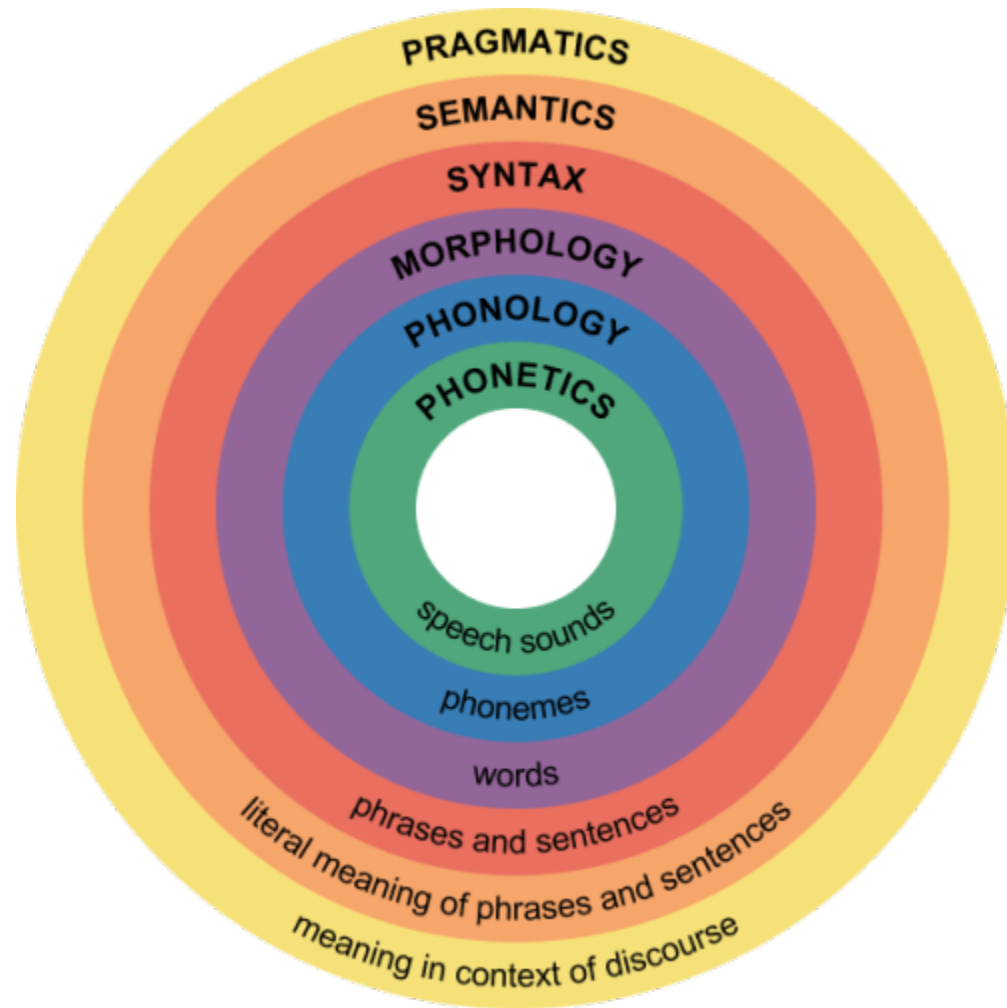
1. It understands you as well as your mother understands you
2. It understands (that) you like your mother
3. It understands you as well as it understands your mother

How many meanings can you recognize in this sentence?

One morning I shot an elephant in my pajamas.

Linguistic components

Sound ->
phoneme ->
morpheme ->
sub-word ->
word ->
phrase ->
sentence ->
paragraph ->
text ->
corpus



<https://pediaa.com/difference-between-semantics-and-pragmatics/>

Definitions of some Linguistic terms

- **Phoneme** a unit of sound that can distinguish one word from another in a particular language, e.g. the word 'cat' has three phonemes: /c/ /a/ /t/.
- **Morpheme** the smallest meaningful lexical item in a language. For example, "un-", "break", and "-able" in the word "unbreakable"
- **Sub-word** just a part of word, like morphemes.
- **Phrase** a group of words or singular word acting as a grammatical unit. For example, adjective phrase "very happy".
- **Corpus** a language resource consisting of a large and structured set of texts, e.g., a set of Wikipedia articles.
- **Semantics** meaning in a literal sense, e.g., "5 pm" means literally a time point.
- **Pragmatics** is the meaning in context, e.g., "5 pm" may mean it is time to go home.

<https://en.wikipedia.org/wiki/>

<http://www.cse.unsw.edu.au/~billw/nlpdict.html>

Parts of Speech

<https://www.sketchengine.eu/penn-treebank-tagset/>

POS Tag	Description	Example
CC	coordinating conjunction	and
CD	cardinal number	1, third
DT	determiner	the
IN	preposition, subordinating conjunction	in, of, like
JJ	adjective	green
NN	noun, singular or mass	table
NNS	noun plural	tables
NP	proper noun, singular	John
NPS	proper noun, plural	Vikings
PP	personal pronoun	I, he, it
VB	verb be, base form	be

Natural Language Processing

NLU

Syntactic parsing

Coreference resolution (“Elon Musk” \Leftrightarrow ”Tesla CEO”)

Semantic parsing (meaning)

Part-of-speech tagging (POS)

Named entity recognition (NER) (“Rome” \Leftrightarrow ”capital of Italy”)

Natural language inference (true, neutral, false)

Relation extraction (“player wins game”)

Text categorization

Sentiment analysis

....

NLU & NLG

Paraphrase

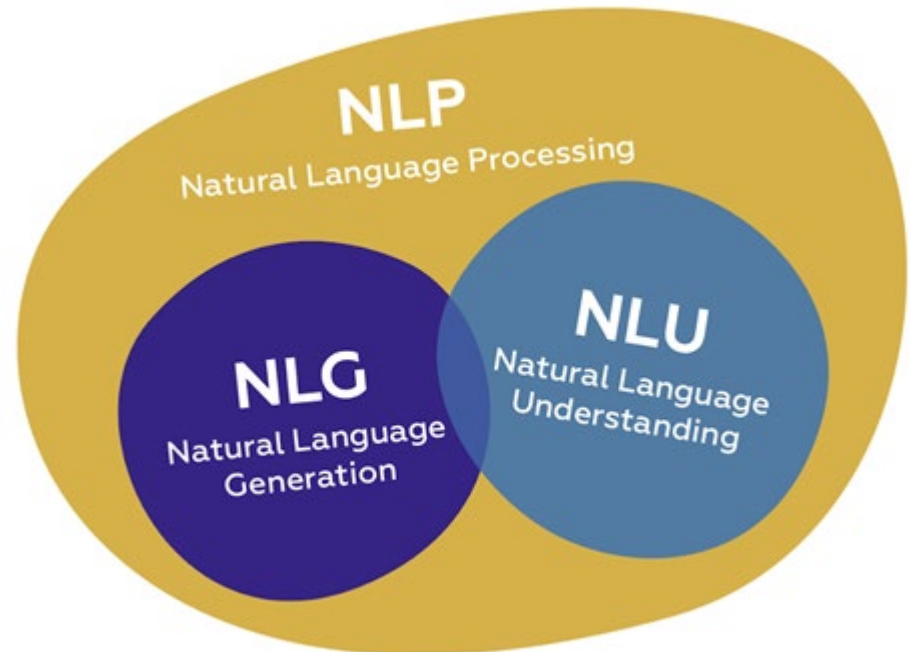
Dialogue agents

Question answering


Text summarization

Machine translation

...



Text classification



US National Library of Medicine
National Institutes of Health

PubMed

Advanced

Display Settings: ☒ Abstract [Send to:](#) ☒

Nature, 2014 Mar 20;507(7492):323-8. doi: 10.1038/nature13145. Epub 2014 Mar 12.

Coupling of angiogenesis and osteogenesis by a specific vessel subtype in bone.

Kusumbe AP¹, Ramasamy SK¹, Adams RH².

Author information

Abstract

The mammalian skeletal system harbours a hierarchical system of mesenchymal stem cells, osteoblasts, that drive bone formation. Osteogenesis is indispensable for the homeostatic renewal of bone as osteoblasts frequently decline in ageing organisms, leading to loss of bone mass and increased fracture risk. Here, we show that blood vessels in bone and osteogenesis are coupled, but relatively little is known about how they interact. We identify a new capillary subtype in the murine skeletal system with distinctive morphological and functional properties. These vessels, specific locations, mediate growth of the bone vasculature and osteogenesis, mediate bone growth, and osteoprogenitors and couple angiogenesis to osteogenesis in bone from aged animals, and pharmacological treatment of aged animals with anti-angiogenic agents.

Comment in

Bone biology: Vascular biology in bone. [View Article](#)

PMID: 24644444



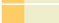
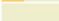
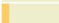
MeSH Terms
[Aging/metabolism](#)
[Aging/pathology](#)
[Animals](#)
[Blood Vessels](#)
[Bone](#)

Sentiment Analysis

Customer Reviews

Speech and Language Processing, 2nd Edition

15 Reviews

5 star:  (8)
4 star:  (3)
3 star:  (3)
2 star:  (0)
1 star:  (1)

Average Customer Review

★★★★☆ (15 customer reviews)

Share your thoughts with other customers

Create your own review

The most helpful favorable review

4 of 4 people found the following review helpful

★★★★★ Great Introductions and reference book

I read the first edition of that book and it is terrific. The second edition is much more adapted to current research. Statistical methods in NLP are more detailed and some syntax-based approaches are presented. My specific interest is in machine translation and dialogue systems. Both chapters are extensively rewritten and much more elaborated. I believe this book is...

[Read the full review >](#)

Published on August 9, 2008 by carheg

> See more [5 star](#), [4 star](#) reviews

Vs.

The most helpful critical review

37 of 37 people found the following review helpful

★★★☆☆ Good description of the problems in the field, but look elsewhere for practical solutions

The authors have the challenge of covering a vast area, and they do a good job of highlighting the hard problems within individual sub-fields, such as machine translation. The availability of an accompanying Web site is a strong plus, as is the extensive bibliography, which also includes links to freely available software and resources.

Now for the...

[Read the full review >](#)

Published on April 2, 2009 by P. Nadkarni

> See more [3 star](#), 2 star, [1 star](#) reviews

Information Extraction



lancet

a Medication Event Extraction System for Clinical Text

Project Home

[Downloads](#)

[Wiki](#)

[Issues](#)

[Source](#)

Summary [People](#)

Project Information

★ Starred by 1 user
[Project feeds](#)

Code license
[GNU GPL v2](#)

Labels
medication, extractor,
lancet, discharge,
summary, i2b2, NLP,
challenge, 2009




Members
[lizuof...@gmail.com](#)

Lancet is a supervised machine-learning system that automatically extracts medication events consisting of medication names and information pertaining to their prescribed use (dosage, mode, frequency, duration and reason) from lists or narrative text in medical discharge summaries.

[Thus, she was transitioned over to a ciprofloxacin 700 mg p.o. b.i.d. regime for a total of 12 days for a presumed urinary tract infection.] narrative


 = medication = dosage = manner = frequency = duration = reason

Information retrieval

 panama papers  

All Images Shopping News Videos More ▾ Search tools


About 88.000.000 results (0,57 seconds)

Datenleak Panama Papers - sueddeutsche.de
 www.sueddeutsche.de/panamapapers ▾
Alle Details zu den Enthüllungen jetzt mit SZ Plus lesen
Bleiben Sie informiert · Alle News zum Thema · Immer aktuell

The Panama Papers · ICIJ
<https://panamapapers.icij.org/> ▾
Politicians, Criminals and the Rogue Industry That Hides Their Cash.

Panama Papers - Wikipedia, the free encyclopedia
https://en.wikipedia.org/wiki/Panama_Papers ▾
The **Panama Papers** are a leaked set of 11.5 million confidential documents that provide detailed information about more than 214,000 offshore companies ...

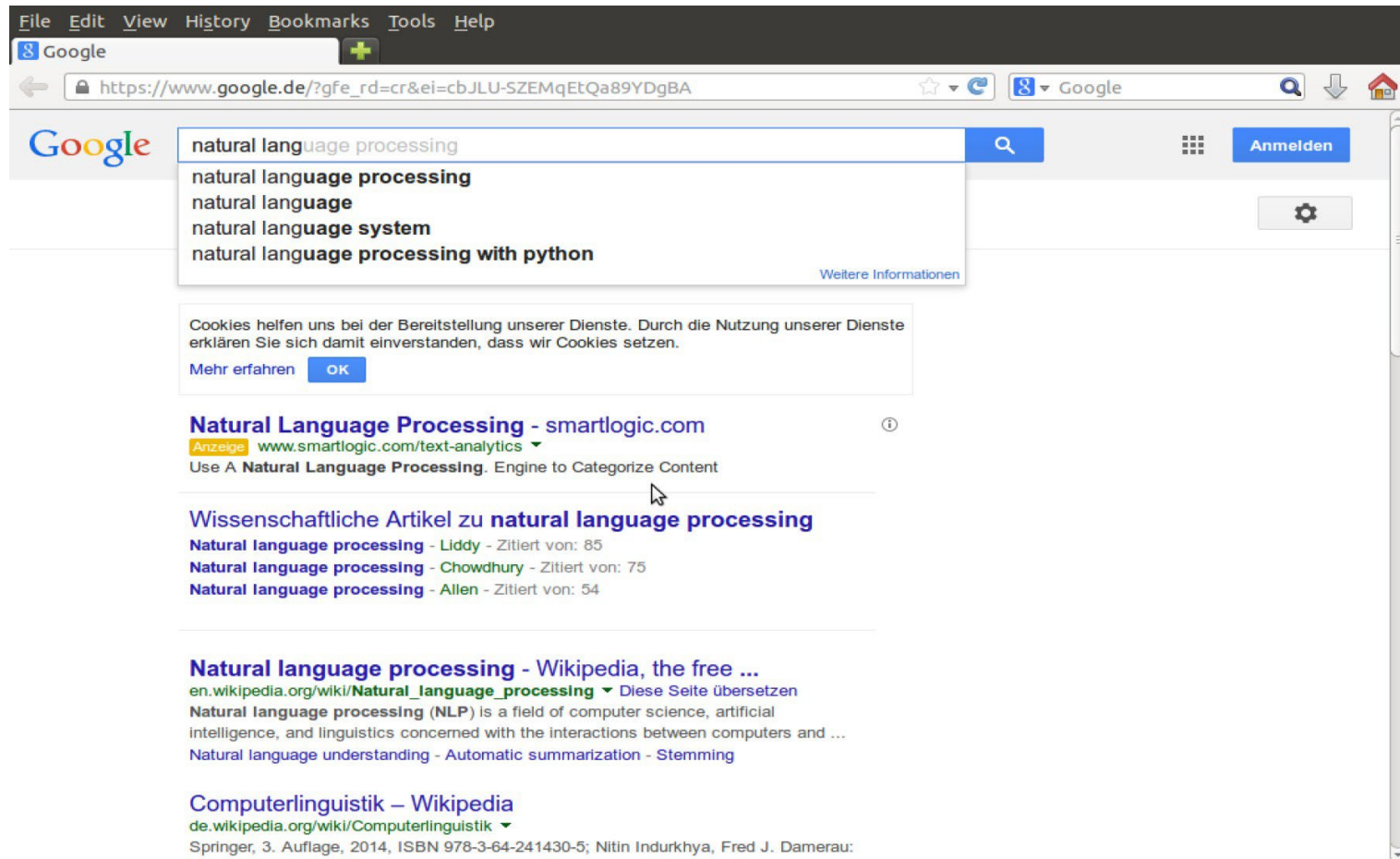
In the news



Panama Papers: Putin rejects corruption allegations - BBC News
BBC News - 2 hours ago
President Putin has denied "any element of corruption" over the **Panama Papers** leaks, ...

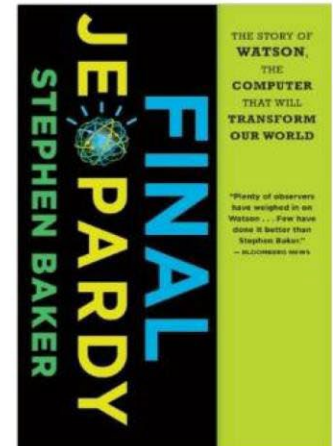
Panama Papers: David Cameron admits profiting from fund
...

Word prediction



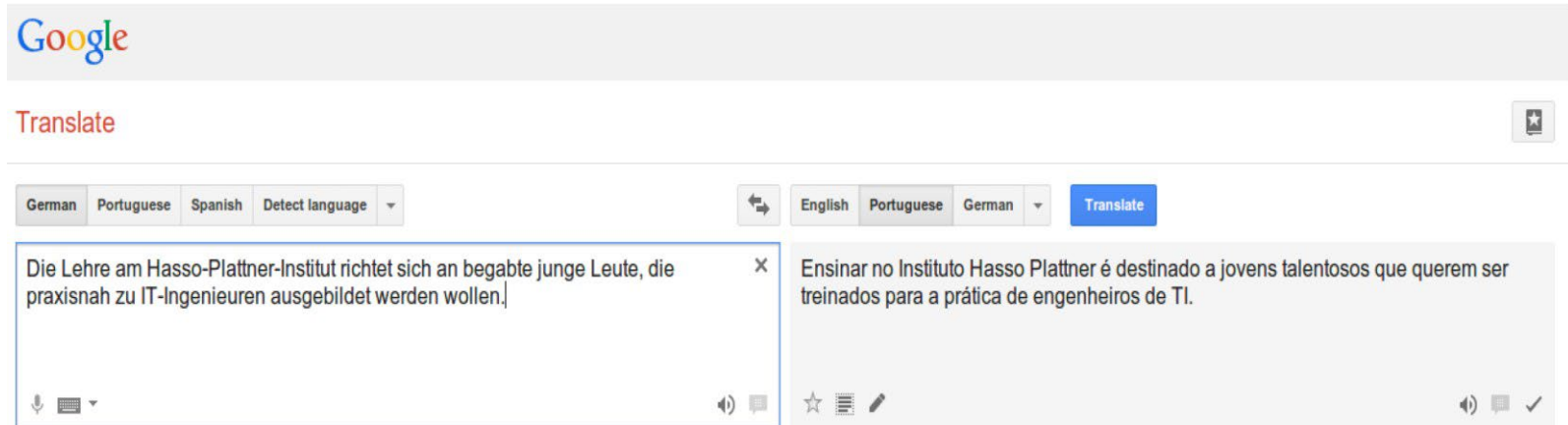
Question Answering

- IBM Watson in Jeopardy



https://www.youtube.com/watch?v=WFR3IOm_xhE

Machine Translation



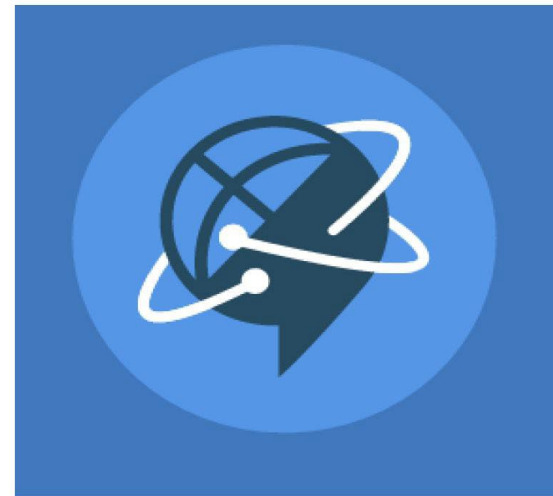
Spoken Dialog System



Siri.
Your wish is
its command.

Siri lets you use your voice to send messages, schedule meetings, place phone calls, and more. Ask Siri to do things just by talking the way you talk. Siri understands what you say, knows what you mean, and even talks back. Siri is so easy to use and does so much, you'll keep finding more and more ways to use it.

IBM Watson Developer Cloud



Summarization

Automatic Text Summarizer

Best Online Summarizing Tool

round from earlier periods. In 1950, Alan Turing published an article titled "Computing Machinery and Intelligence" which proposed what is now called the Turing test as a criterion of intelligence[clarification needed].

The Georgetown experiment in 1954 involved fully automatic translation of more than sixty Russian sentences into English. The authors claimed that within three or five years, machine translation would be a solved problem.[2] However, real progress was much slower, and after the ALPAC report in 1966, which found that ten-year-long research had failed to fulfill the expectations, funding for machine translation was

Clear

Summarize

Some notably successful natural language processing systems developed in the 1960s were SHRDLU, a natural language system working in restricted blocks worlds with restricted vocabularies, and ELIZA, a simulation of a Rogerian psychotherapist, written by Joseph Weizenbaum between 1964 and 1966.

Up to the 1980s, most natural language processing systems were based on complex sets of hand-written rules.

Starting in the late 1980s, however, there was a revolution in natural language processing with the introduction of machine

Recap of Machine Learning

1. What are the similarities and differences between logistic regression and neural networks?
2. Given the function $f(x) = x^3$, what is the value of the slope at $x = 1$?
3. What is Stochastic Gradient Descent, and how is it functionally different to Batch Gradient Descent?
4. How does regularisation reduce overfitting? Explain this given the cost function with regularisation.

$$L_1 = (wx + b - y)^2 + \lambda|w|$$

$$L_2 = (wx + b - y)^2 + \lambda w^2$$

<https://towardsdatascience.com/intuitions-on-l1-and-l2-regularisation-235f2db4c261>

5. Share your impressions of using the tensorflow playground. What have you learned, what was interesting?