

Natural Language Processing

SoSe 2017



Introduction to Natural Language Processing

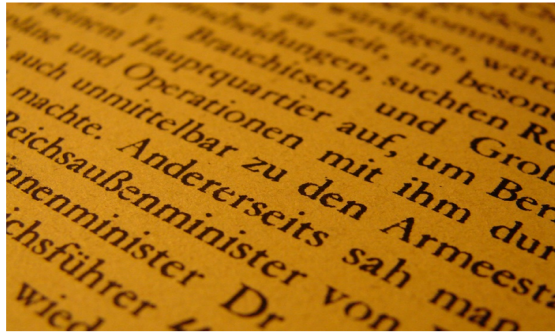
Dr. Mariana Neves

April 24th, 2017

Outline

- Introduction to NLP
- NLP Applications
- NLP Techniques
- Linguistic Knowledge
- Challenges
- NLP course

Natural Language



日本語で

ふゆ せかいかくち さまざまなお祝いが 行 われる時期で
 す。ほんのいくつか 例 を挙げるだけでも、ハナカ、クリス
 マス、クワンザ、 新年 などさまざまなお祝いが あります。
 各文化によってその 祝い方 はさまざまですが、ほとん
 どのお祝 いにはごちそうが 欠かせません。

(<http://expertenough.com/2392/german-language-hacks>)

(http://www.transparent.com/learn-japanese/articles/dec_99.html)

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) {
    // ...
}

```

- append(CharSequence)
- append(char c)
- append(CharSequence)
- format(String format, Object... args)
- format(Locale l, String format, Object... args)
- printf(String format, Object... args)

```

def add5(x):
    return x+5

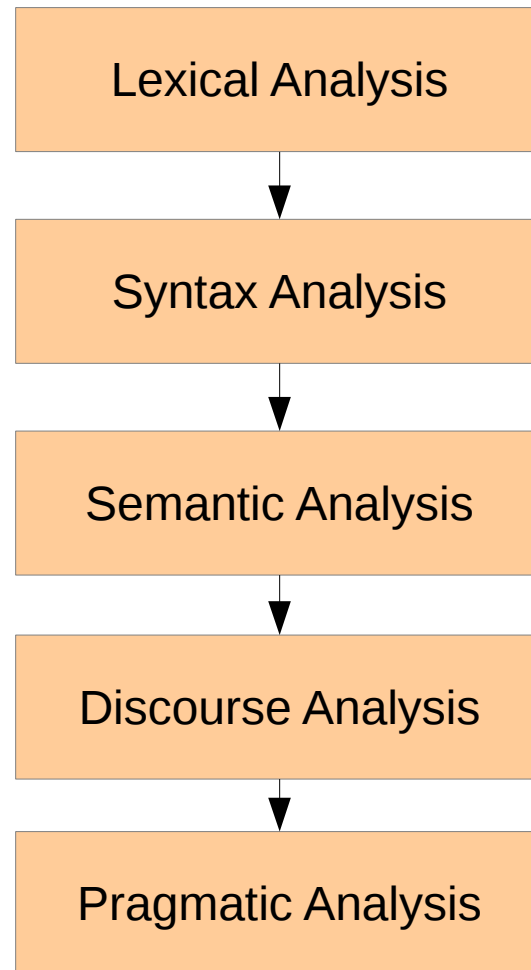
def dotwrite(ast):
    nodename = getNodeName()
    label=symbol.sym_name.get(int(ast[0]),ast[0])
    print '  %s [label="%s" % (nodename,label),
    if isinstance(ast[1], str):
        if ast[1].strip():
            print ' = %s';' % ast[1]
        else:
            print ''
    else:
        print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

```

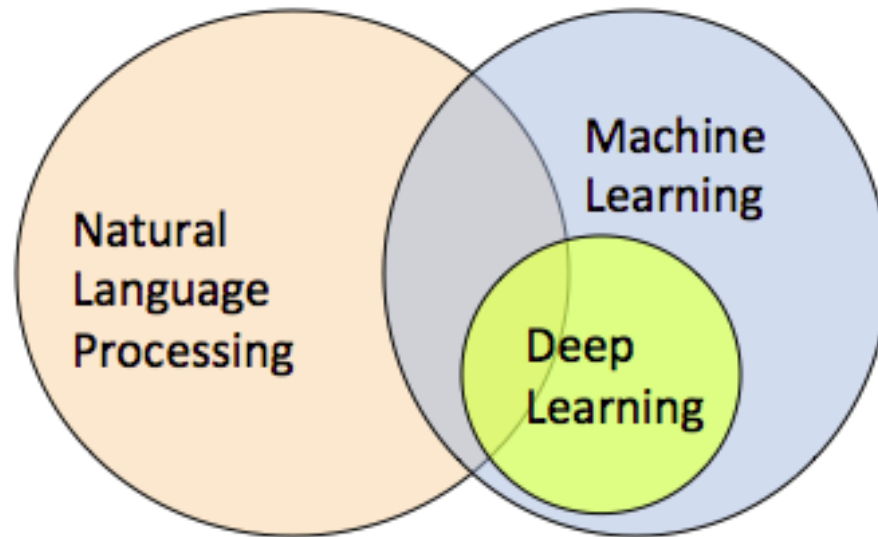
(<http://noobite.com/learn-programming-start-with-python/>)

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

Natural Language Processing



Natural Language Processing

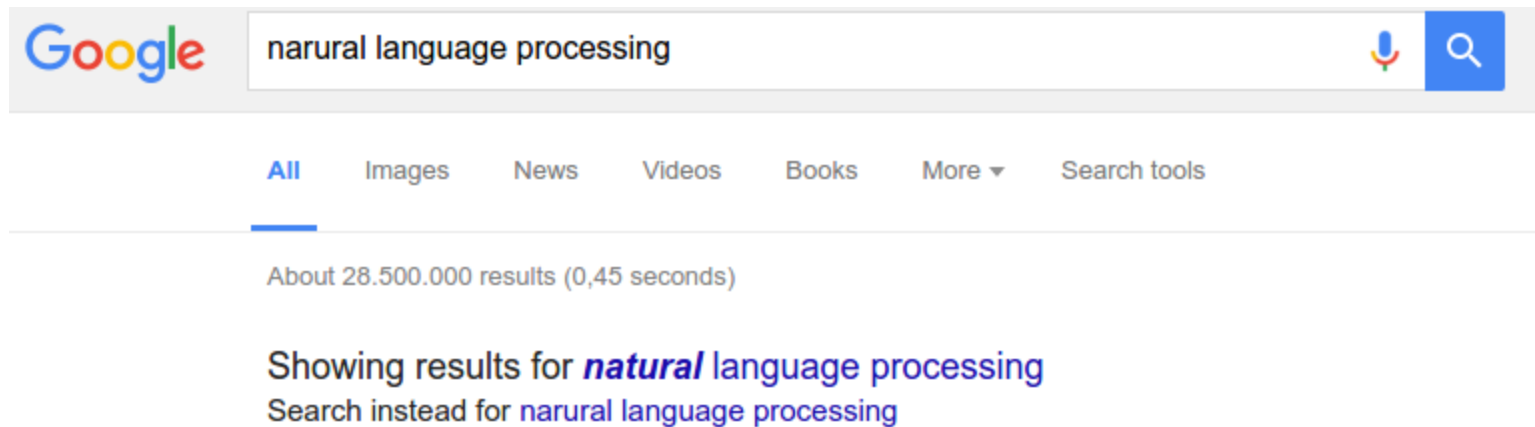


Outline

- Introduction to NLP
- **NLP Applications**
- NLP Techniques
- Linguistic Knowledge
- Challenges
- NLP course

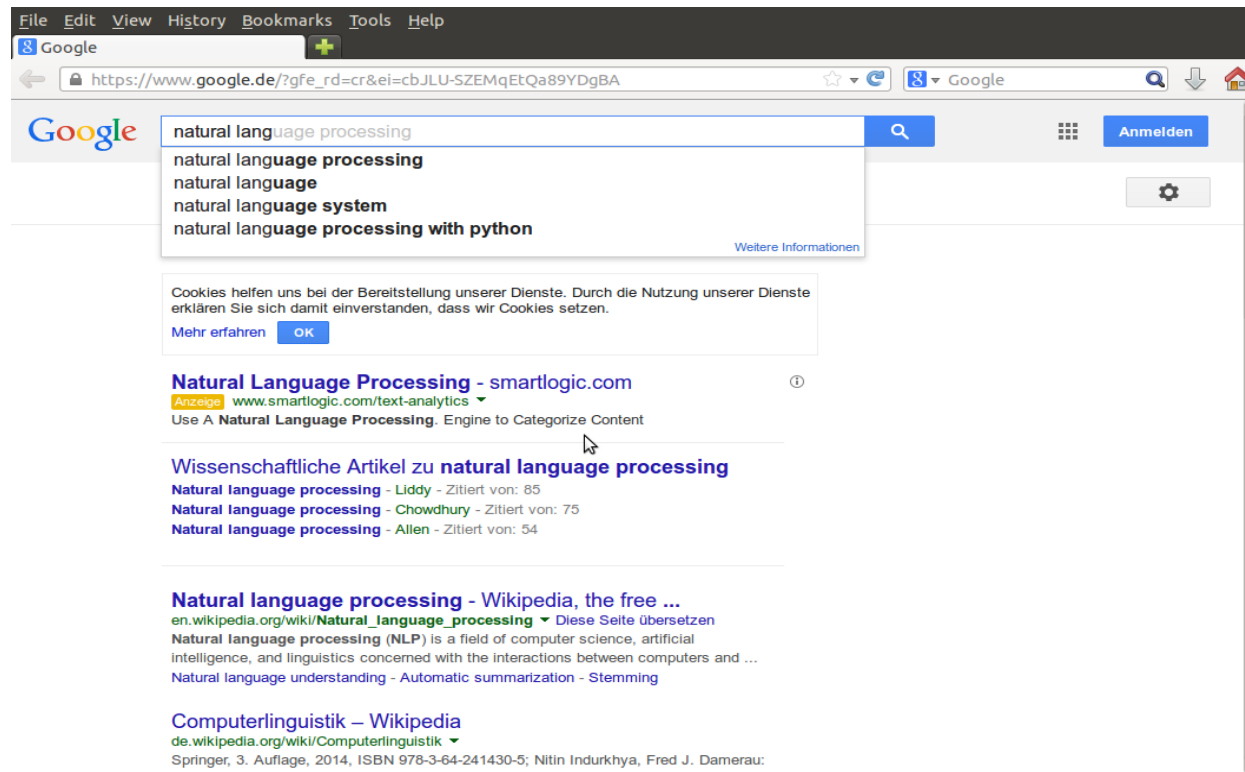
Spell and Grammar Checking

- Checking spelling and grammar
- Suggesting alternatives for the errors



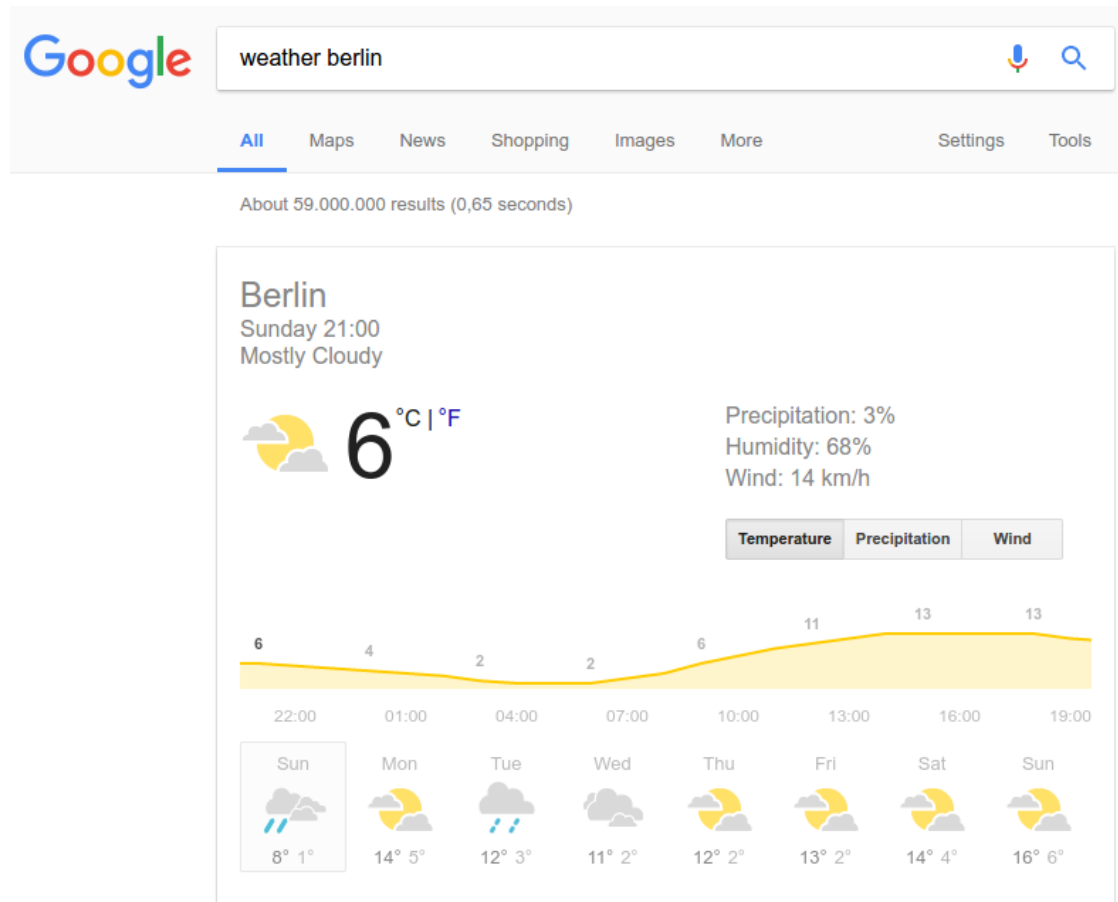
Word Prediction

- Predicting the most probable next word to be typed by the user



Information Retrieval

- Finding relevant information to the user's query



Information Retrieval

- Finding relevant information to the user's query

The screenshot displays the Olelo search interface. At the top, there is a search bar with the text "What are the diseases caused by the zika virus?". Below the search bar, the results are organized into two main panels. The left panel, titled "List of Diseases", shows a list of diseases with "MICROCEPHALY" highlighted in blue. The right panel, titled "Abstracts of 15/41 documents", shows a list of abstracts. The first abstract is titled "Zika virus: history of a newly emerging arbovirus." and includes the author "DR Smith, N Wilan", the journal "Lancet Infect Dis, 2016, 16(7)", and the PMID "27282424". The abstract text describes the history of Zika virus, its transmission, and its effects on humans. At the bottom of the abstract panel, there are two buttons: "CREATE SUMMARY" and "EXPORT AS BIOC".

Text Categorization

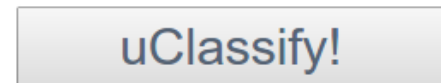
- Assigning one (or more) pre-defined category to a text



Classify

Classify method: text url

Enter url to download and classify with:



Remove html

1. Sports (92.8 %)
2. Entertainment (4.8 %)
3. Men (0.7 %)

[Show all classifications >>](#)

Summarization

- Generating a short summary from one or more documents, sometimes based on a given query



This is a sentence summary of <http://hpi.de/en/news/jahrgaenge/2015/des...>

Summary processing at low priority, upgrade to **BOOST**

Design Thinking Week: Students Improve the Daily Life Experience for People with Illiteracies

On the occasion of the World Literacy Day on September 8 more than 40 young innovators applied their Design Thinking skills in order to make life easier for these people.

Here, the focus was especially on the possibilities of using digital technologies and computers to better the daily obstacles in life of the people concerned.

Under the guidance of the D-School's coaches the teams researched, developed and prototyped - and could present many versatile solutions in the end: e.g. one of the groups came up with an idea for a software program that lets internet browsers read texts, functions and links out loud so that people with reading problems can still use news sites or social networks like Facebook.

Summarization

- Generating a short summary from one or more documents, sometimes based on a given query

Olelo Search

Abstracts of 15/41 documents

Your search: Microcephaly

Zika virus: history of a newly emerging arbovirus.

DR Smith, N Wikan
Lancet Infect Dis, 2016, 16(7)
PMID: 27282424

Zika virus was originally identified in a sentinel rhesus monkey in the Zika Forest of Uganda in 1947. The virus is a member of the family Flaviviridae, genus Flavivirus, and is transmitted to humans by Aedes species mosquitoes. The first report of Zika virus outside Africa and Asia was in 2007 when the virus was associated with a small outbreak in Yap State, part of the Federated States of Micronesia. Since then, Zika virus infections have been reported around the world, including in southeast Asia; French Polynesia and other islands in the Pacific Ocean; and parts of South, Central, and North America. Symptomatic infection in human beings normally results in a mild and self-limiting febril...

CREATE SUMMARY EXPORT AS BIOC

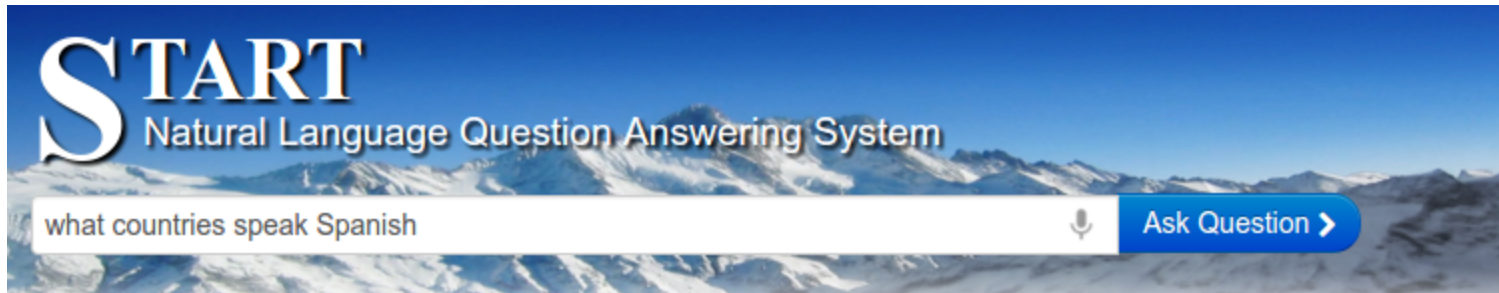
Summary of 15 documents

It is transmitted by the Aedes genus of mosquitoes. Definite cases had laboratory evidence of Zika virus infection; highly probable cases presented specific neuroimaging findings, and negative laboratory results for other congenital infections; moderately probable cases had specific imaging findings but other infections could not be ruled out; somewhat probable cases had imaging findings, but these were not reported in detail by the local teams; all other newborn babies were classified as discarded cases. We further show that ZIKV infection, but not WNV infection, impairs cell cycle progression of neural stem cells. Historically, ZIKV infection was characterized by a self-limiting, mild disease, but recent outbreaks have been associated with severe clinical complications, including Guillain-Barré syndrome and microcephaly, which are atypical of other flavivirus infections. The widespread outbreak and accelerating increase in the number of cases in Puerto Rico warrants intensified vector control and personal protective behaviors to prevent new infections, particularly among pregnant women.

TRANSLATE CORRESPONDING DOCUMENTS

Question answering

- Answering questions with a short answer



==> what countries speak Spanish

The language Spanish is spoken in Argentina, Aruba, Belize, Bolivia, Brazil, Canada, Cayman Islands, Chile, Colombia, Costa Rica, Cuba, Curacao, Dominican Republic, Ecuador, El Salvador, Equatorial Guinea, Falkland Islands (Islas Malvinas), Gibraltar, Guatemala, Honduras, Mexico, Nicaragua, Panama, Paraguay, Peru, Puerto Rico, Saint Martin, Sint Maarten, Spain, Switzerland, Trinidad and Tobago, United States, Uruguay, Venezuela, and Virgin Islands.

The language Castilian Spanish is spoken in Spain.

Question Answering

Olelo

What are the diseases caused by the zika virus?

List of Diseases

- DENGUE
- RIFT VALLEY FEVER
- MICROCEPHALY
- ZIKA VIRUS INFECTION
- YELLOW FEVER
- INFECTION
- BITES AND STINGS
- ARTHRALGIA
- EMERGENCIES
- NERVOUS SYSTEM DISEASES
- MALARIA
- CONJUNCTIVITIS

SELECT A TERM TO SHOW AMOUNT OF MATCHING DOCUMENTS

How to treat Guillain-Barre syndrome?

A child with Guillain-Barre syndrome treated with intravenous immune globulin (IVIg) developed neutropenia (absolute neutrophil count = 390), which resolved 3 days after completion of the therapy. This report deals with an elderly lady with Guillain-Barre syndrome (GBS), who presented with features of unusually severe hyponatraemia. Gangliosides are abundantly expressed in the nervous system, and deregulated expression or activity of Gangliosides is associated with the progression of various disorders, including lysosomal storage diseases, Guillain-Barre syndrome and Alzheimer disease. As there is no specific drug for GBS, several drugs targeting the humoral and cellular components of the immune response have been used to treat EAN in the endeavour to find new treatment alternatives for GBS. At that time, the mean improvement was 0.9 (SD 1.3) in the 121 PE-group patients, 0.8 (1.3) in the 130 IVIg-group patients, and 1.1 (1.4) in the 128 patients who received both treatments (intention-to-treat analysis).

TRANSLATE CORRESPONDING DOCUMENTS

Question answering

- IBM Watson in Jeopardy



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

Information Extraction

- Extracting important concepts from texts and assigning them to slot in a certain template

Angela Merkel



Merkel at the EPP Summit, March 2016

Chancellor of Germany

Incumbent

Assumed office
22 November 2005

President [Horst Köhler](#)
[Christian Wulff](#)
[Joachim Gauck](#)

Deputy [Franz Müntefering](#)
[Frank-Walter Steinmeier](#)
[Guido Westerwelle](#)
[Philipp Rösler](#)
[Sigmar Gabriel](#)

Preceded by [Gerhard Schröder](#)

Leader of the Christian Democratic Union

Incumbent

Assumed office
10 April 2000

Preceded by [Wolfgang Schäuble](#)

Minister for the Environment

In office

17 November 1994 – 26 October 1998

Chancellor [Helmut Kohl](#)

Preceded by [Klaus Töpfer](#)

Succeeded by [Jürgen Trittin](#)

Minister for Women and Youth

In office

18 January 1991 – 17 November 1994

Chancellor [Helmut Kohl](#)

Preceded by [Ursula Lehr](#)

Succeeded by [Claudia Nolte](#)

Personal details

Born [Angela Dorothea Kasner](#)
17 July 1954 (age 61)
[Hamburg, West Germany](#)

Political party [Democratic Awakening](#) (1989–1990)
[Christian Democratic Union](#) (1990–present)

Spouse(s) [Ulrich Merkel](#) (1977–1982)
[Joachim Sauer](#) (1998–present)

Alma mater [Leipzig University](#)

Religion [Lutheranism](#) (within [Evangelical Church](#))

Signature



Information Extraction

- Includes named-entity recognition

Angela Merkel

From Wikipedia, the free encyclopedia

"Merkel" redirects here. For other uses, see [Merkel \(disambiguation\)](#).

Angela Dorothea Merkel (English /ˈæŋɡələ ˈmɜːrkəl/^[a]; née **Kasner**; born 17 July 1954) is a German politician who is currently [Chancellor of Germany](#). She is also the leader of the [Christian Democratic Union](#) (CDU). Merkel has been described at various times as the *de facto* leader of the [European Union](#), the most powerful woman in the world, and the world's [second most powerful person](#).

A former [research scientist](#) with a doctorate in [physical chemistry](#), Merkel entered politics in the wake of the [Revolutions of 1989](#), and briefly served as a deputy spokesperson for the first democratically elected [East German Government](#) headed by [Lothar de Maizière](#) in 1990. Following [German reunification](#) in 1990, Merkel was elected to the [Bundestag](#) for the state of [Mecklenburg-Vorpommern](#), and has been reelected ever since. Merkel was appointed as the [Minister for Women and Youth](#) in the [federal government](#) under [Chancellor Helmut Kohl](#) in 1991, and became the [Minister for the Environment](#) in 1994. After her party lost the [federal election in 1998](#), Merkel was elected [Secretary-General](#) of the CDU before becoming the party's first woman leader two years later in the aftermath of a [donations scandal](#) that toppled [Wolfgang Schäuble](#).

Following the [2005 federal election](#), Merkel was appointed Germany's first woman Chancellor at the head of a [grand coalition](#) consisting of the CDU, its Bavarian sister party, the [Christian Social Union](#) (CSU), and the [Social Democratic Party of Germany](#) (SPD). In the [2009 federal election](#), the CDU obtained the largest share of the vote and Merkel was able to form a coalition government with the support of the [Free Democratic Party](#) (FDP).^[9] At the [2013 federal election](#), Merkel's CDU won a landslide victory with 41.5% of the vote and formed a second grand coalition with the SPD, after the FDP lost all of its representation in the [Bundestag](#).^[10]

In 2007, Merkel was [President of the European Council](#) and chaired the [G8](#), the second woman to do so. Merkel played a central role in the negotiation of the [Treaty of Lisbon](#) and the [Berlin Declaration](#). One of Merkel's consistent priorities has been to strengthen transatlantic economic relations. Merkel played a crucial role in managing the [financial crisis](#) at the European and international level, and she has been referred to as "the decider." In domestic policy, [health care reform](#), problems concerning future [energy development](#) and more recently her government's approach to the ongoing [migrant crisis](#) have been major issues during her Chancellorship.^[11] On 26 March 2014, Merkel became the longest-serving incumbent [head of government](#) in the [European Union](#) and she is currently the [senior G7](#) leader. On 20 November 2016, Merkel announced she would seek re-election to a fourth term.^[12]

(https://en.wikipedia.org/wiki/Angela_Merkel)

Information Extraction

Gynecol Oncol. 2007 Dec;107(3):518-25. Epub 2007 Oct 25.

Combination therapy with pegylated liposomal doxorubicin and carboplatin in gynecologic malignancies: a prospective phase II study of the Arbeitsgemeinschaft Gynäkologische Onkologie Studiengruppe Ovarialkarzinom (AGO-OVAR) and Kommission Uterus (AGO-K-Ut).

du Bois A¹, Pfisterer J, Burchardi N, Loibl S, Huober J, Wimberger B, Burger A, Stöbke A, Iseliach C, Kölbl H; Arbeitsgemeinschaft Gynäkologische Onkologie Studiengruppe Ovarialkarzinom; Kommission U

⊕ Author information

Abstract

OBJECTIVE: A multicenter non-randomized phase doxorubicin (PLD) in combination with carboplatin

METHODS: One hundred forty women with recurr sarcomas (n=11), or recurrent platinum-sensitive every 28 days.

RESULTS: Hematological toxicities with NCI-CTC febrile neutropenia in 2% of 652 cycles. Grade 3/4 (9%), palmar-plantar erythrodysesthesia (7%), an carboplatin. Seventy-four percent of all non-progn evaluable for response): ovarian cancer (n=54) 6 cancer (n=26) 12%. Median progression-free surv 12.6) for endometrial cancer. Median overall survi

CONCLUSIONS: The combination of PLD and carl Efficacy was low in cervical/vaginal cancer, but pr recurrent platinum-sensitive ovarian cancer and is comparing PLD/carboplatin with paclitaxel/carbop

PMID: 17910981 DOI: [10.1016/j.ygyno.2007.08.008](https://doi.org/10.1016/j.ygyno.2007.08.008)

Evidence / Table

ID	[[25345_dubois2007]]			
Kapitel	Vaginalkarzinom			
	Endometriumkarzinom Ovarialkarzinom Zervixkarzinom			
Zuordnung	4_cht_pall_vagina 4_cht_pall_endo 4_cht_pall_ovar 4_cht_pall_zervix	PDF		
Typ	NCT			
Intervention(en)	Cth	Pall	Pegyliertes liposomales Doxorubicin +	40 mg/m ² + AUC 6 Alle 4 Wochen

(<https://www.ncbi.nlm.nih.gov/pubmed/17910981>)

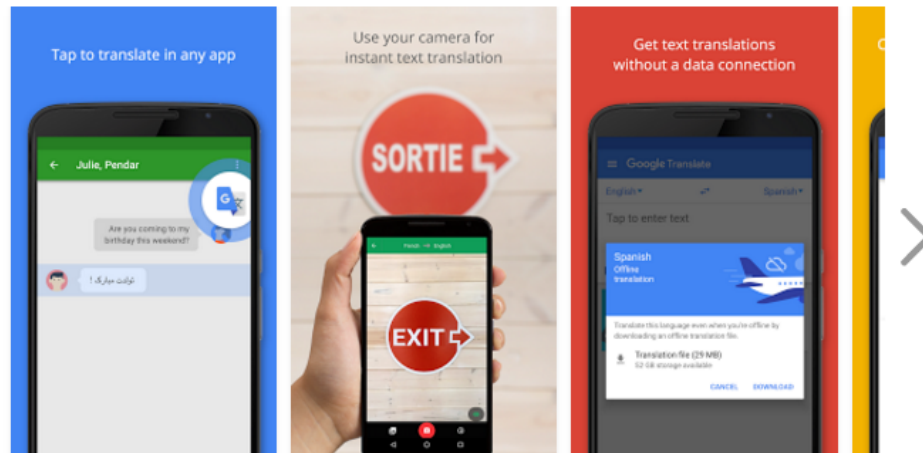
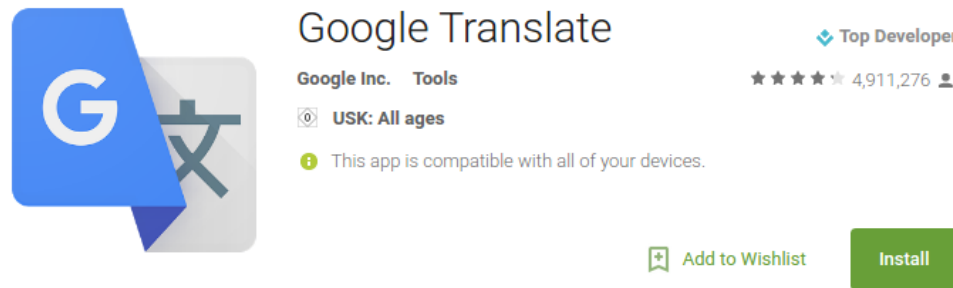
Machine Translation

- Translating a text from one language to another

The image is a screenshot of a news article from Le Monde. At the top, there is a navigation bar with the text "AT THE MOMENT" in a red box, followed by links for "Presidential Election 2017", "Syria", "Donald Trump", "Case Fillon", "French, French", "The decoders", and "Big Browser". The main headline is "Presidency 2017: Macron and Le Pen qualified, Fillon and Hamon call to vote Macron" with "LIVE" in a small font to the right. Below the headline is a sub-headline: "According to a new estimate Ipsos-Sopra Steria for 'Le Monde', the candidate of En marche! Would obtain 23.9% of the vote and that of the National Front 21.7%. They are ahead of François Fillon (20%) and Jean-Luc Mélenchon (19.2%)." The main image shows two portraits side-by-side: Emmanuel Macron on the left and Marine Le Pen on the right. A red "LIVE" button is visible in the bottom left corner of the image area.

Machine Translation

- Translating a text from one language to another



Sentiment Analysis

- Identifying sentiments and opinions stated in a text

Customer Reviews Speech and Language Processing, 2nd Edition



The most helpful favorable review	vs.	The most helpful critical review
<p>4 of 4 people found the following review helpful</p> <p>★★★★★ 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...</p> <p>Read the full review > Published on August 9, 2008 by carheg</p> <p>> See more 5 star, 4 star reviews</p>		<p>37 of 37 people found the following review helpful</p> <p>★★★☆☆ 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.</p> <p>Now for the... Read the full review > Published on April 2, 2009 by P. Nadkarni</p> <p>> See more 3 star, 2 star, 1 star reviews</p>

Optical Character Recognition

- Recognizing printed or handwritten texts and converting them to computer-readable texts



Speech recognition

- Recognizing a spoken language and transforming it into a text



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.

Speech synthesis

- Producing a spoken language from a text



Spoken dialog systems

- Running a dialog between the user and the 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.

Level of difficulties

- Easy (mostly solved)
 - Spell and grammar checking
 - Some text categorization tasks
 - Some named-entity recognition tasks

Level of difficulties

- Intermediate (good progress)
 - Information retrieval
 - Sentiment analysis
 - Machine translation
 - Information extraction

Level of difficulties

- Difficult (still hard)
 - Question answering
 - Summarization
 - Dialog systems

Outline

- Introduction to NLP
- NLP Applications
- **NLP Techniques**
- Linguistic Knowledge
- Challenges
- NLP course

Section splitting

• Splitting a text into sections

Eur Radiol
DOI 10.1007/s00330-014-3135-8

BREAST

Correlation between three-dimensional ultrasound features and pathological prognostic factors in breast cancer

Jun Jiang · Yi-qing Chen · Yi-zhan Xu · Ming-E Chen · Yun-kai Zhu · Wen-bin Guan · Xiao-jin Wang

Received: 13 November 2013 / Rev. recd.: 30 January 2014 / Accepted: 17 February 2014
© European Society of Radiology 2014

Abstract
Objectives To investigate the correlation of three-dimensional (3D) ultrasound features with prognostic factors in invasive ductal carcinoma.
Methods Surgical resection specimens of 85 invasive ductal carcinomas of 85 women who had undergone 3D ultrasound were included. Morphology features and vascularization perfusion on 3D ultrasound were evaluated. Pathologic prognostic factors, including tumour size, histological grade, lymph node status, oestrogen and progesterone receptor status (ER, PR), c-erbB-2 and p53 expression, and microvessel density (MVD) were determined. Correlations of 3D ultrasound features and prognostic factors were analysed.
Results The retraction pattern in the coronal plane had a significant value as an independent predictor of a small tumour size ($P=0.014$), a lower histological grade ($P=0.009$) and positive ER or PR expression status ($P=0.001$, 0.044). The retraction pattern with a hypercholeic ring only existed in low-grade and ER-positive tumours. The presence of the hypercholeic ring strengthened the ability of the retraction pattern to predict a good prognosis of breast cancer. The increased intra-tumour vascularization index (VI, the mean

tumour vascularity) reflected a higher histological grade ($P=0.025$) and had a positive correlation with MVD ($r=0.530$, $P=0.001$).
Conclusions The retraction pattern and histogram indices of VI provided by 3D ultrasound may be useful in predicting prognostic information about breast cancer.
Key Points
• Three-dimensional ultrasound can potentially provide prognostic evaluation of breast cancer.
• The retraction pattern and hypercholeic ring in the coronal plane suggest good prognosis.
• The increased intra-tumour vascularization index reflects a higher histological grade.
• The intra-tumour vascularization index is positively correlated with microvessel density.

Keywords Breast · Neoplasms · Ultrasound · Three-dimensional · Prognostic factors

Introduction

The three strongest prognostic factors in invasive breast cancer are widely accepted to be the size of tumour, histological grade and lymph node stage. The larger tumour size (>2 cm), high nuclear grade, and lymph node-positive status usually predict the aggressive biological behaviour with a high recurrence rate and a low survival rate. In addition, the tumour size and lymph node status greatly influence the choice of operative procedure and the decision to administer neoadjuvant chemotherapy [1, 2].

Biological markers such as oestrogen receptors (ER), progesterone receptors (PR), human epidermal growth factor receptor 2 (c-erbB-2) and the p53 index can also be used for prediction of medical treatment response and patient prognosis. The presence of ER and PR in breast cancer always

determines the application of antihormonal therapy and usually indicates a good prognosis. Expression of c-erbB-2 or the p53 index is a powerful and independent prognostic factor for lymph node metastasis and tumour infiltration [1, 3]. Microvessel density (MVD) is the current reference standard in the characterization of tumour angiogenesis and has been shown to be associated with lymph node growth, invasion, metastasis and disease-specific survival [4].
Three-dimensional (3D) ultrasound can afford additional information such as morphology features on the coronal plane and a global appearance of the mass vascularity, which cannot be achieved with conventional ultrasound. Therefore, it has been increasingly considered as an important imaging modality for evaluating primary breast cancer. However, so far, 3D ultrasound has been used mainly to differentiate benign and malignant lesions; no reports address correlations between the 3D ultrasound features and prognostic factors [5–7]. We therefore investigated possible correlation between the 3D ultrasound characteristics of invasive ductal carcinoma with pathologic prognostic factors to determine whether 3D ultrasound could be useful in the non-invasive prognostic evaluation of breast cancer.

Materials and methods

Patients

This retrospective study was approved by the ethical standards of the institutional ethics committee, and informed consent was obtained from all patients.
From September 2011 to May 2013, 85 patients with 85 lesions, pathologically proven to be invasive ductal carcinoma, were included in this study. The exclusion criteria were pregnancy or lactation, administration of preoperative chemotherapies or adjuvant chemotherapies. Patients with a breast mass larger than 3.0 cm were also excluded because more than one 3D volume acquisition was necessary to include the whole lesion plus 3 mm surrounding the breast lesion. All patients were female and aged 26 to 90 years (mean age, 56.3 years).

Ultrasound examination

All ultrasound images were obtained with one type of system (GE Voluson E8 Expert, Zipf, Austria) by two radiologists with 7–12 years of experience in breast ultrasound. An 11 L-D linear transducer with a frequency of 5–12 MHz was used for 2D ultrasound, and an RSPr-16-D dedicated volume transducer with a frequency of 6–12 MHz was used for 3D ultrasound.

Ultrasound examination was performed with patients in the supine position with elevated arms. Once the breast lesion was

detected and the region of interest had been identified, the volume box was superimposed and set to include the entire display screen so as to cover the lesion and maximum amount of normal surrounding tissue. The sweep angle was adjusted to 15–29° according to the size of the breast lesion. Then the ultrasound probe was held still with enough jelly to contact the skin gently. The volume mode was switched on and the 3D ultrasound volume was generated by the automatic rotation of the mechanical transducer. When the first ultrasound examination was finished, the power Doppler mode was added for the second examination and the fixed preinstalled power Doppler settings used were 0.3 kHz pulse repetition frequency, “low 1” wall motion filter, “2.0 gain and high frequency. The first examination for 3D grayscale imaging took 10–20 s and the second, for 3D power Doppler imaging, took 25–45 s, depending on the size of the tumour. Then the total acquisition time for 3D ultrasound was about 1–2 min. The entire examination was saved in DICOM format and stored on the hard disk for further analysis.

Image analysis

The 3D ultrasound images were reviewed for this analysis by another two radiologists with 8–10 years of experience in breast ultrasound and characterized by consensus. In addition, the radiologists had not performed the data acquisition and were blinded to the patients’ clinical and mammographic findings.

The ultrasound image was opened by using the 4D View software. Firstly, the tomographic ultrasound imaging (TUI) was used for a slice by slice documentation in the coronal plane. Then, the volume contrast imaging (VCI) and the surface render mode were added for better observation of the lesion and the surrounding tissue. All the slices were carefully observed to identify the presence of the retraction pattern in the surrounding tissue and the margin of the lesion. The retraction pattern was defined as the hypercholeic straight lines that radiated perpendicularly from the surface of the solid nodule, producing a stellate pattern [8, 9] (Fig. 1). The presence of the retraction pattern was further divided into with or without a hypercholeic ring, which was displayed as an echogenic halo ring between the mass and the surrounding tissue in the coronal plane (Fig. 2a).

The 3D power Doppler imaging analyses were performed using a virtual organ computer-aided analysis (VOCAL)-imaging program (GE, Zipf, Austria), which could automatically calculate the histogram indices of vascularization index (VI), flow index (FI) and vascularization flow index (VFI). VI represents the vessels in the defined volume by measuring the number of colour voxels in the region of interest, i.e. the mean tumour vascularity; FI represents the average intensity of flow by measuring the mean colour value in the colour voxels, i.e. the mean blood flow volume; VFI represents both

Eur Radiol

Eur Radiol

regression modelling techniques to identify the most significant and independent 3D image findings. A P value less than 0.05 was considered statistically significant.

Results

Prognostic factors

In the current study group, the surgical specimens revealed 75 lesions with pure invasive ductal carcinoma and the remaining 10 lesions with invasive ductal carcinoma with DCIS components. The mean percentage of the DCIS components in the lesion was 8.10±4.93% (range, 2–20%).
The size of 85 lesions ranged from 5 to 30 mm, and the mean size was 19.92 mm (SD=7.56 mm). Of the 85 tumours, 47 (55.3%) were equal to or smaller than 2 cm and 38 (44.7%) were larger than 2 cm. According to the Elston–Ellis grading system, there were 58 (68.2%) grade II tumours and 27 (31.8%) grade III. Lymph node metastasis was present in 30 (35.3%) patients. There were 58 (68.2%) ER-positive, 54 (63.5%) PR-positive, 70 (82.4%) c-erbB-2-positive and 42 (49.4%) p53-positive tumours.

Correlation between MVD and prognostic factors

Significantly higher MVD was observed in the larger size group ($P<0.01$) and higher grade group ($P<0.05$). There were no significant associations between MVD and other pathological factors ($P>0.05$) (Table 1).

Correlation between morphological features and prognostic factors

Of the 85 breast lesions, 57 (67.1%) showed the retraction pattern in the coronal plane of 3D ultrasound. Of these 57 lesions, 17 (29.8%) showed the retraction pattern with a hypercholeic ring and 40 (70.2%) were without the hypercholeic ring.

The tumour size, histological grade, ER and PR status all showed significant associations with the presence of the retraction pattern ($P<0.01$) (Table 2). Tumours with the retraction pattern were significantly more likely to be small in size, low grade, ER-positive and PR-positive (Fig. 3). Moreover, the retraction pattern with a hypercholeic ring, which presented as intricately mixed fibrous tissues and infiltrating carcinoma cells on pathological specimens, only existed in low-grade and ER-positive tumours (Fig. 2). The odds ratios of tumour size, tumour grade, and ER and PR status for patients with the retraction pattern and a hypercholeic ring versus no retraction pattern were all higher than those with the retraction pattern without a hypercholeic ring versus no retraction pattern (Table 3). The presence of the hypercholeic ring strengthened

Table 1 Association between MVD and prognostic factors

Prognostic factor	N	Mean	SD	P value
Tumour size (cm)				
≤2	47	19.30	5.25	
>2	38	25.60	7.60	0.007
Tumour grade				
II	58	19.83	5.55	
III	27	25.83	8.02	0.023
Lymph node				
Negative	55	21.31	6.70	
Positive	30	22.08	7.34	0.946
ER				
Negative	27	23.27	8.36	
Positive	58	20.93	5.14	0.931
PR				
Negative	31	25.00	8.59	
Positive	54	19.82	5.09	0.092
c-erbB-2				
Negative	15	21.50	9.57	
Positive	70	21.55	6.65	0.788
p53				
Negative	43	23.13	7.04	
Positive	42	19.63	6.20	0.083

the ability of the retraction pattern to predict these good prognoses. However, the lymph node status and the expression of c-erbB-2 and p53 showed no statistically significant correlation with the retraction pattern ($P>0.05$).

As for MVD, however, no significant correlation was found between MVD and the presence of the retraction pattern on 3D ultrasound ($P=0.05$).

Correlation between vascularization perfusion and prognostic factors

For intra-tumour regions, the mean VI, FI and VFI of 85 lesions were 6.84 (range, 0.02–21.61), 37.72 (range, 21.81–53.32) and 2.64 (range, 0.04–9.11), respectively. For shells with a thickness of 3 mm surrounding the breast lesion, the VI, FI and VFI were 7.31 (range, 0.14–25.13), 38.72 (range, 23.27–56.90) and 2.88 (range, 0.04–11.08), respectively.

Compared with the small tumours, the tumour foci with a diameter greater than 2 cm were more likely to show a higher inVI, inFI, inVFI, out3mmVI and out3mmVFI. The tumours with a high grade or lymph node metastasis had a higher inVI, inVFI, out3mmVI and out3mmVFI than the tumours with low grade or lymph node-negative status. ER-negative tumours had a higher inFI than ER-positive tumours and the tumours with negative expression of PR had a higher inVI, inVFI and out3mmVFI than PR-positive tumours (Table 4).

Published online: 12 April 2014



Sentence splitting

- Splitting a text into sentences

11 Sentences (= "T-" or "Terminable" units *only* if independent clauses are punctuated as separate sentences, e.g. "I came and he went"-->"I came. And he went.")

Average 23.55 words (SD=12.10)

OBJECTIVES: To investigate the correlation of three-dimensional (3D) ultrasound features with prognostic factors in invasive ductal carcinoma.

METHODS: Surgical resection specimens of 85 invasive ductal carcinomas of 85 women who had undergone 3D ultrasound were included.

Morphology features and vascularization perfusion on 3D ultrasound were evaluated.

Pathologic prognostic factors, including tumour size, histological grade, lymph node status, oestrogen and progesterone receptor status (ER, PR), c erbB-2 and p53 expression, and microvessel density (MVD) were determined.

Correlations of 3D ultrasound features and prognostic factors were analysed.

RESULTS: The retraction pattern in the coronal plane had a significant value as an independent predictor of a small tumour size (P #8201;= 0.014), a lower histological grade (P #8201;= 0.009) and positive ER or PR expression status (P #8201;= 0.001, 0.044).

The retraction pattern with a hyperechoic ring only existed in low-grade and ER-positive tumours.

The presence of the hyperechoic ring strengthened the ability of the retraction pattern to predict a good prognosis of breast cancer.

The increased intra-tumour vascularization index (VI, the mean tumour vascularity) reflected a higher histological grade (P #8201;= 0.025) and had a positive correlation with MVD (r #8201;= 0.530, P #8201;= 0.001).

CONCLUSIONS: The retraction pattern and histogram indices of VI provided by 3D ultrasound may be useful in predicting prognostic information about breast cancer.

KEY POINTS: • Three-dimensional ultrasound can potentially provide prognostic evaluation of breast cancer. • The retraction pattern and hyperechoic ring in the coronal plane suggest good prognosis. • The increased intra-tumour vascularization index reflects a higher histological grade. • The intra-tumour vascularization index is positively correlated with microvessel density.

Part-of-speech tagging

- Assigning a syntactic tag to each word in a sentence

Stanford Parser

Please enter a sentence to be parsed:

Surgical resection specimens of 85 invasive ductal carcinomas of 85 women who had undergone 3D ultrasound were included.

Language: [Sample Sentence](#)

Your query

Surgical resection specimens of 85 invasive ductal carcinomas of 85 women who had undergone 3D ultrasound were included.

Tagging

Surgical/NNP resection/NN specimens/NNS of/IN 85/CD invasive/JJ ductal/JJ carcinomas/NNS of/IN 85/CD women/NNS who/WP had/VBD undergone/VBN 3D/CD ultrasound/NN were/VBD included/VBN ./.

Parsing

- Building the syntactic tree of a sentence

Parse

```

(ROOT
  (S
    (NP
      (NP (NNP Surgical) (NN resection) (NNS specimens))
      (PP (IN of)
        (NP
          (NP (CD 85) (JJ invasive) (JJ ductal) (NNS carcinomas))
          (PP (IN of)
            (NP
              (NP (CD 85) (NNS women))
              (SBAR
                (WHNP (WP who))
                (S
                  (VP (VBD had)
                    (VP (VBN undergone)
                      (NP (CD 3D) (NN ultrasound))))))))))
            (VP (VBD were)
              (VP (VBN included)))
          (. .)))

```

Parsing

- Building the syntactic tree of a sentence

Typed dependencies

```
nn(specimens-3, Surgical-1)
nn(specimens-3, resection-2)
nsubjpass(included-18, specimens-3)
prep(specimens-3, of-4)
num(carcinomas-8, 85-5)
amod(carcinomas-8, invasive-6)
amod(carcinomas-8, ductal-7)
pobj(of-4, carcinomas-8)
prep(carcinomas-8, of-9)
num(women-11, 85-10)
pobj(of-9, women-11)
nsubj(undergone-14, who-12)
aux(undergone-14, had-13)
rcmod(women-11, undergone-14)
num(ultrasound-16, 3D-15)
dobj(undergone-14, ultrasound-16)
auxpass(included-18, were-17)
root(ROOT-0, included-18)
```

Named-entity recognition

- Identifying pre-defined entity types in a sentence

The screenshot displays the b2cas Annotate interface. On the left, a 'HIGHLIGHT' sidebar lists categories: Anatomy, Disorders, Chemicals, Genes and Proteins, Cellular Components, Molecular Functions, Biological Processes, and Ambiguous. The main text area contains a paragraph about Duchenne muscular dystrophy (DMD) with various entities highlighted in colored boxes. Below the text, a 'Load text' button and an 'Export' button are visible. At the bottom, a 'Concept Tree' shows a hierarchical view of the annotated entities, including categories like Anatomy (12), Disorders (4), Chemicals (2), Genes and Proteins (11), Cellular Components (3), Molecular Functions (1), and Biological Processes (9).

Text from screenshot:

In **Duchenne muscular dystrophy (DMD)**, the **infiltration** of **skeletal muscle** by immune **cells** aggravates disease, yet the precise mechanisms behind these **inflammatory responses** remain poorly understood. Chemotactic cytokines, or chemokines, are considered essential recruiters of **inflammatory cells** to the **tissues**. We assayed chemokine and chemokine **receptor expression** in **DMD muscle** biopsies (n = 9, average age 7 years) using immunohistochemistry, immunofluorescence, and in situ **hybridization**. **CXCL1**, **CXCL2**, **CXCL3**, **CXCL8**, and **CXCL11**, absent from normal **muscle fibers**, were induced in **DMD** myofibers. **CXCL11**, **CXCL12**, and the **ligand**-receptor couple **CCL2**-**CCR2** were upregulated on the **blood vessel endothelium** of **DMD** patients. **CD68** (+) **macrophages** expressed high levels of **CXCL8**, **CCL2**, and **CCL5**. Our data suggest a possible beneficial role for **CXCR1/2/4 ligands** in managing **muscle fiber** damage control and **tissue** regeneration. Upregulation of **endothelial chemokine receptors** and **CXCL8**, **CCL2**, and **CCL5** expression by cytotoxic **macrophages** may regulate myofiber necrosis.

Concept Tree:

- Anatomy (12)
 - Disorders (4)
 - DMD (1)
 - Duchenne muscular dystrophy (1)
 - infiltration (1)
 - inflammatory responses (1)
 - Chemicals (2)
 - Genes and Proteins (11)
 - Cellular Components (3)
 - Molecular Functions (1)
 - Biological Processes (9)

Word sense disambiguation

- Figuring out the exact meaning of a word or entity

Noun 1. tie - neckwear consisting of a long narrow piece of material worn (mostly by men) under a collar and tied in knot at the front; "he stood in front of the mirror tightening his necktie"; "he wore a vest and tie"

[necktie](#)

[bola](#), [bola tie](#), [bolo](#), [bolo tie](#) - a cord fastened around the neck with an ornamental clasp and worn as a necktie

[bow tie](#), [bow-tie](#), [bowtie](#) - a man's tie that ties in a bow

[four-in-hand](#) - a long necktie that is tied in a slipknot with one end hanging in front of the other

[neckwear](#) - articles of clothing worn about the neck

[old school tie](#) - necktie indicating the school the wearer attended

[string tie](#) - a very narrow necktie usually tied in a bow

[Windsor tie](#) - a wide necktie worn in a loose bow

2. tie - a social or business relationship; "a valuable financial affiliation"; "he was sorry he had to sever his ties with other members of the team"; "many close associations with England"

[affiliation](#), [tie-up](#), [association](#)

[relationship](#) - a state involving mutual dealings between people or parties or countries



3. tie - equality of score in a contest

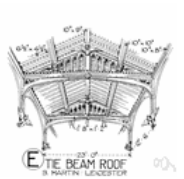
[equivalence](#), [par](#), [equality](#), [equation](#) - a state of being essentially equal or equivalent; equally balanced; "on a par with the best"

[deuce](#) - a tie in tennis or table tennis that requires winning two successive points to win the game

4. tie - a horizontal beam used to prevent two other structural members from spreading apart or separating; "he nailed the rafters together with a tie beam"

[tie beam](#)

[beam](#) - long thick piece of wood or metal or concrete, etc., used in construction



Word sense disambiguation

Analysis with definitions(s)

Bill Gates has developed an interest/[readiness to give attention] in language technology and yesterday acquired a 10 % interest/[a share (in a company, business, etc.)] in Torbjörn Lager 's sense disambiguation technology . Lager will retain a 90 % interest/[a share (in a company, business, etc.)] in the new company , which will be based in Göteborg , Sweden . Last year 's drop in interest/[money paid for the use of money] rates will probably be good for the company . Finally , although all this may sound like an arcane maneuver of little interest/[quality of causing attention to be given] outside Wall Street , it would set off an economical earthquake .

These are the six senses of the noun *interest* according to the LDOCE:

Sense	Definition
1	readiness to give attention
2	quality of causing attention to be given
3	activity, subject, etc., which one gives time and attention to
4	advantage, advancement, or favour
5	a share (in a company, business, etc.)
6	money paid for the use of money

Semantic role labeling

- Extracting subject-predicate-object triples from a sentence



Semantic Role Labeling Demo

Input Text:

They had brandy in the library .

[Click For General Explanation of Argument Labels](#)

Output:

	<input type="checkbox"/> SRL	<input type="checkbox"/> Nom	<input type="checkbox"/> Preposition	<input type="checkbox"/>
They	owner [A0]			
had	V: have.03			
brandy	possession [A1]		Governor	
in			Locationin:1(1)	
the	location [AM-LOC]			
library			Object	
.				

Outline

- Introduction to NLP
- NLP Applications
- NLP Techniques
- Linguistic Knowledge
- Challenges
- NLP course

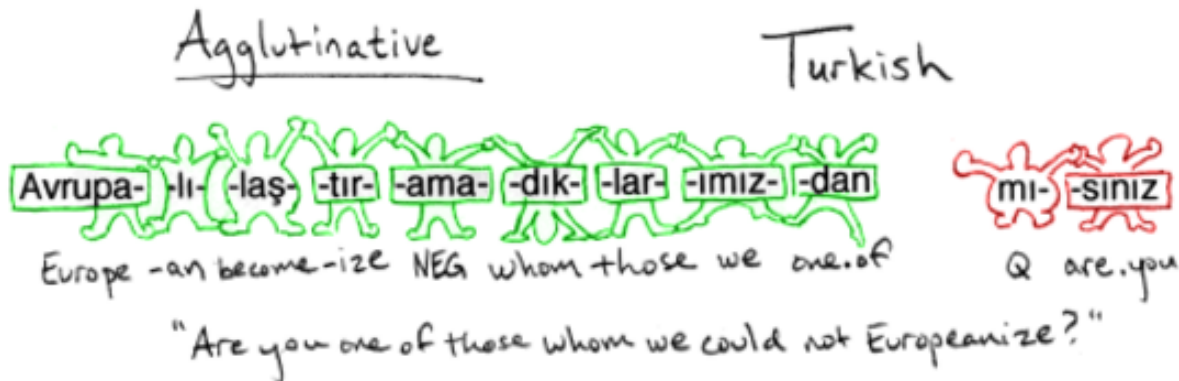
Phonetics and phonology

- The study of linguistic sounds and their relations to words

Das Funkalphabet - German Phonetic Spelling Code compared to the international ICAO/NATO code Listen to AUDIO for this chart! (below)		
Germany*	Phonetic Guide	ICAO/NATO**
A wie Anton	AHN-tone	Alfa/Alpha
Ä wie Ärger	AIR-gehr	(1)
B wie Berta	BARE-tuh	Bravo
C wie Cäsar	SAY-zar	Charlie
Ch wie Charlotte	shar-LOT-tuh	(1)
D wie Dora	DORE-uh	Delta
E wie Emil	ay-MEAL	Echo
F wie Friedrich	FREED-reech	Foxtrot
G wie Gustav	GOOS-tahf	Golf
H wie Heinrich	HINE-reech	Hotel
I wie Ida	EED-uh	India/Indigo
J wie Julius	YUL-ee-oos	Juliet
K wie Kaufmann	KOWF-mann	Kilo
L wie Ludwig	LOOD-vig	Lima
AUDIO 1 > Listen to mp3 for A-L		
M wie Martha	MAR-tuh	Mike
N wie Nordpol	NORT-pole	November
O wie Otto	AHT-toe	Oscar
Ö wie Ökonom (2)	UEH-ko-nome	(1)
P wie Paula	POW-luh	Papa
Q wie Quelle	KVEL-uh	Quebec
R wie Richard	REE-shart	Romeo
S wie Siegfried (3)	SEEG-freed	Sierra
Sch wie Schule	SHOO-luh	(1)
ß (Eszett)	ES-TSET	(1)
T wie Theodor	TAY-oh-dore	Tango
U wie Ulrich	OOL-reech	Uniform
Ü wie Übermut	UEH-ber-moot	(1)
V wie Viktor	VICK-tor	Victor
W wie Wilhelm	VIL-helm	Whiskey
X wie Xanthippe	KSAN-tipp-uh	X-Ray
Y wie Ypsilon	IPP-see-lohn	Yankee
Z wie Zeppelin	TSEP-puh-leen	Zulu

Morphology

- The study of internal structures of words and how they can be modified
- Parsing complex words into their components



Syntax

- The study of the structural relationships between words in a sentence

Parse

```
(ROOT
  (S
    (NP
      (NP (NNP Surgical) (NN resection) (NNS specimens))
      (PP (IN of)
        (NP
          (NP (CD 85) (JJ invasive) (JJ ductal) (NNS carcinomas))
          (PP (IN of)
            (NP
              (NP (CD 85) (NNS women))
              (SBAR
                (WHNP (WP who))
                (S
                  (VP (VBD had)
                    (VP (VBN undergone)
                      (NP (CD 3D) (NN ultrasound))))))))))
          (VP (VBD were)
            (VP (VBN included)))
          (. .)))
```

Semantics

- The study of the meaning of words, and how these combine to form the meanings of sentences
 - Synonymy: fall & autumn
 - Hypernymy & hyponymy (is a): animal & dog
 - Meronymy (part of): finger & hand
 - Homonymy: fall (verb & season)
 - Antonymy: big & small

Pragmatics

- Social use of language
- The study of how language is used to accomplish goals, and the influence of context on meaning.
- Understanding the aspects of a language which depends on situation and world knowledge.

“Give me the salt!”

or

“Could you please give me the salt?”

Discourse

- The study of linguistic units larger than a single statement

John reads a book. He borrowed it from his friend.

Berlin (/bɛərˈlɪn/, German: [bɛʁˈliːn] (ⓘ listen)) is the capital of Germany, and one of the 16 states of Germany. With a population of 3.5 million people,^[4] Berlin is Germany's largest city. It is the second most populous city proper and the seventh most populous urban area in the European Union.^[5] Located in northeastern Germany on the banks of River Spree, it is the center of the Berlin-Brandenburg Metropolitan Region, which has about 6 million residents from over 180 nations.^{[6][7][8][9]} Due to its location in the European Plain, Berlin is influenced by a temperate seasonal climate. Around one third of the city's area is composed of forests, parks, gardens, rivers and lakes.^[10]

Outline

- Introduction to NLP
- NLP Applications
- NLP Techniques
- Linguistic Knowledge
- **Challenges**
- NLP course

Paraphrasing

- Different words/sentences express the same meaning
 - Season of the year
 - Fall
 - Autumn
 - Book delivery time
 - When will my book arrive?
 - When will I receive my book?

Ambiguity

- One word/sentence can have different meanings
 - Fall
 - The third season of the year
 - Moving down towards the ground or towards a lower position
 - The door is open.
 - Expressing a fact
 - A request to close the door

Phonetics and Phonology



Communication tip:

Phonological ambiguities or Give peas a chance!

One of my favourite ways to have fun with communication are phonological ambiguities.

Phonological ambiguities are two or more words which sound the same and have different meanings.



Language can contain ambiguities - and more than one way to compose a set of sounds into words.

So listen to yourself: It is always good to notice a spoken sentence often contains many words which are (sometimes not)

intended to be heard.

English examples:

- there - their
- here - hear
- plane - plain
- Hamburger (Citizens of Hamburg) - hamburger (burger, food)
- sea - see
- Friday - fry day
- weekend - weak end
- ice cream - I scream.
- new direction - nude erection
- new day - nude, eh?
- I don't know! - I don't - no!
- but - butt
- Wait - Weight
- psychotherapist - psycho the rapist
- You're unconscious now... - Your unconscious now...
- Your students... - You're students...
- Two - too - to

German examples:

- Du hast Gewehre. (You have got guns.) - Du hasst Gewehre. (You hate guns.)
- Lehrer (teacher) - leerer (emptier)

Syntax and ambiguity

- I saw the man with a telescope.
 - Who had the telescope?



Semantics

- The astronomer loves the **star**.
 - Star in the sky
 - Celebrity



(<http://www.businessnewsdaily.com/2023-celebrity-hiring.html>)

(<http://en.wikipedia.org/wiki/Star#/media/File:Starsinthesky.jpg>)

Discourse analysis

- Alice understands that you like your mother, but **she** ...
 - Does **she** refer to Alice or your mother?

Outline

- Introduction to NLP
- NLP Applications
- NLP Techniques
- Linguistic Knowledge
- Challenges
- **NLP course**

NLP Course

- Home page:
 - <https://hpi.de/plattner/teaching/summer-term-2017/natural-language-processing.html>
- Lecture
 - Monday 11:00-12:30
 - HS3
 - 3 credit points

Grading

- 60% Project
- 40% Final exam (written)

- You have to pass both of them

Program

(Program is subject to change)

Session	Date	Topic
1	April 24, 2017	Introduction to Natural Language Processing
2	May 1, 2017	(no lecture - Maifeiertag)
3	May 8, 2017	Morphology
4	May 15, 2017	Words and Language Model
5	May 22, 2017	Part-of-Speech Tagging and Syntactic Parsing
6	May 29, 2017	Lexical Semantics
7	June 5, 2017	(no lecture - Pfingstmontag)
8	June 12, 2017	Discourse analysis
9	June 19, 2017	Project: mid-term presentation (15 minutes for each team)
10	June 26, 2017	Sentiment Analysis
11	July 3, 2017	Named-entity Recognition and Information Extraction
12	July 10, 2017	Question Answering and Summarization
13	July 17, 2017	Project: final presentation (15 minutes for each team)
14	July 24, 2017	Final exam

Project

- Development of a NLP application
 - Information Retrieval
 - Information Extraction
 - Text Summarization
 - Question Answering
 - Sentiment Analysis
 - Machine Translation
 - Etc..

Project

- The application should include following components:
 - Syntactic parsing (e.g., POS tagging)
 - Semantics (e.g., NER)
 - Discourse analysis

Project

- Any NLP or ML libraries
 - Stanford Core NLP
 - NLTK
 - spaCy
 - Apache OpenNLP
 - GATE
 - SAP HANA (contact me)
 - R
 - Weka
 - TensorFlow
 - etc.

Project

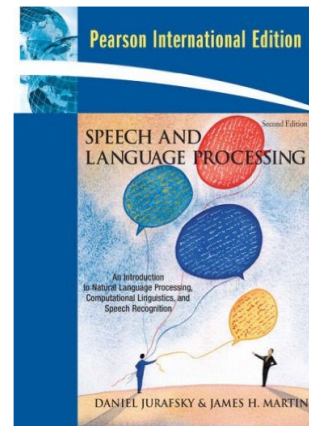
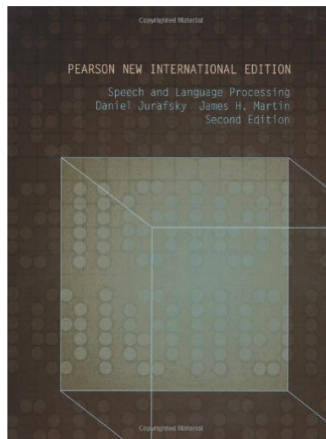
- Any language
 - English, German, etc.
(Check available NLP tools)
- Any text collections
 - Social media, Web pages, publications, Wikipedia, etc.
 - Benchmarks or newly created datasets
- Any domain

Project

- Teams (2-3 students)
- Send me an email with your proposal (until May 5th, 2017)
- Updates (presentations) on the progress of the project
 - Mid-term and final presentation
 - Also considered for grading (commitment)

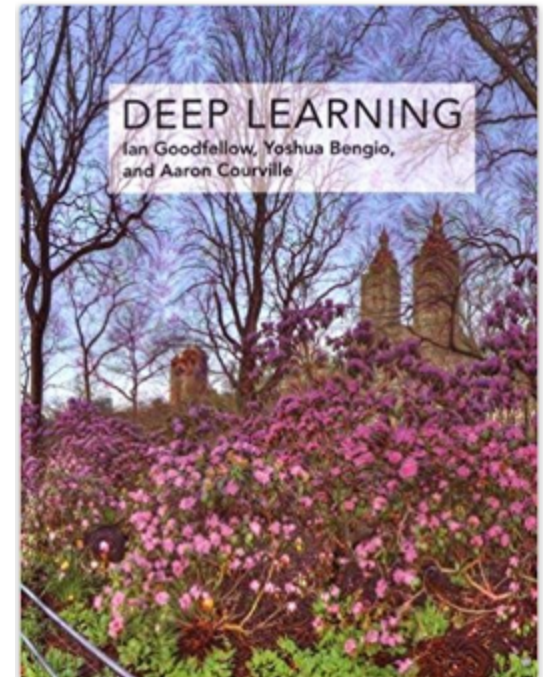
Course book

- Speech and Language Processing
 - Daniel Jurafsky and James H. Martin
 - New (3rd) edition: <https://web.stanford.edu/~jurafsky/slp3/>



Deep learning book

- Deep Learning (Adaptive Computation and Machine Learning)
 - Ian Goodfellow, Yoshua Bengio, Aaron Courville
 - The Mit Press, 2017
 - <http://www.deeplearningbook.org>
 - NLP in section 12.4



Journal and conferences

- Journal
 - Computational Linguistics
- Conferences
 - ACL: Association for Computational Linguistics
 - NAACL: North American Chapter
 - EACL: European Chapter
 - HLT: Human Language Technology
 - EMNLP: Empirical Methods on Natural Language Processing
 - CoLing: Computational Linguistics
 - LREC: Language Resources and Evaluation

NLP Course

- Contact
 - Mariana.Neves@hpi.uni-potsdam.de
 - Room V-0.01 (Villa) - appointment under request

- We have a student position for NLP at the EPIC chair!