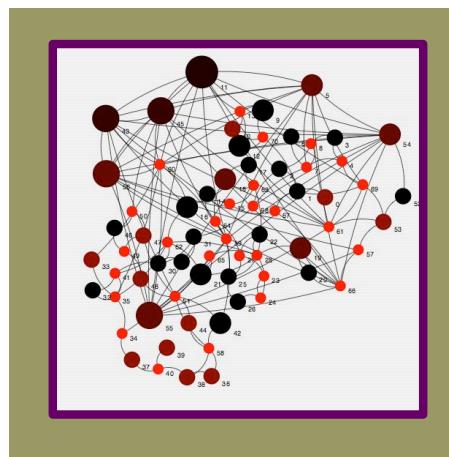
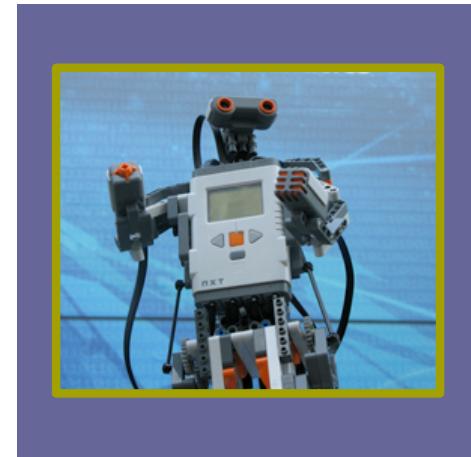
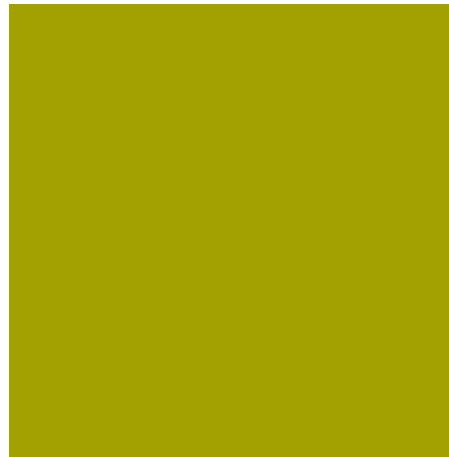




COMP40020

Human Language Technologies

Evolution of Language
Technologies
January 2019



Prof. Julie Berndsen
School of Computer Science

Julie.Berndsen@ucd.ie

Contents:

- Some literature
- Some highlights from the past
- Some current applications and directions

Aim:

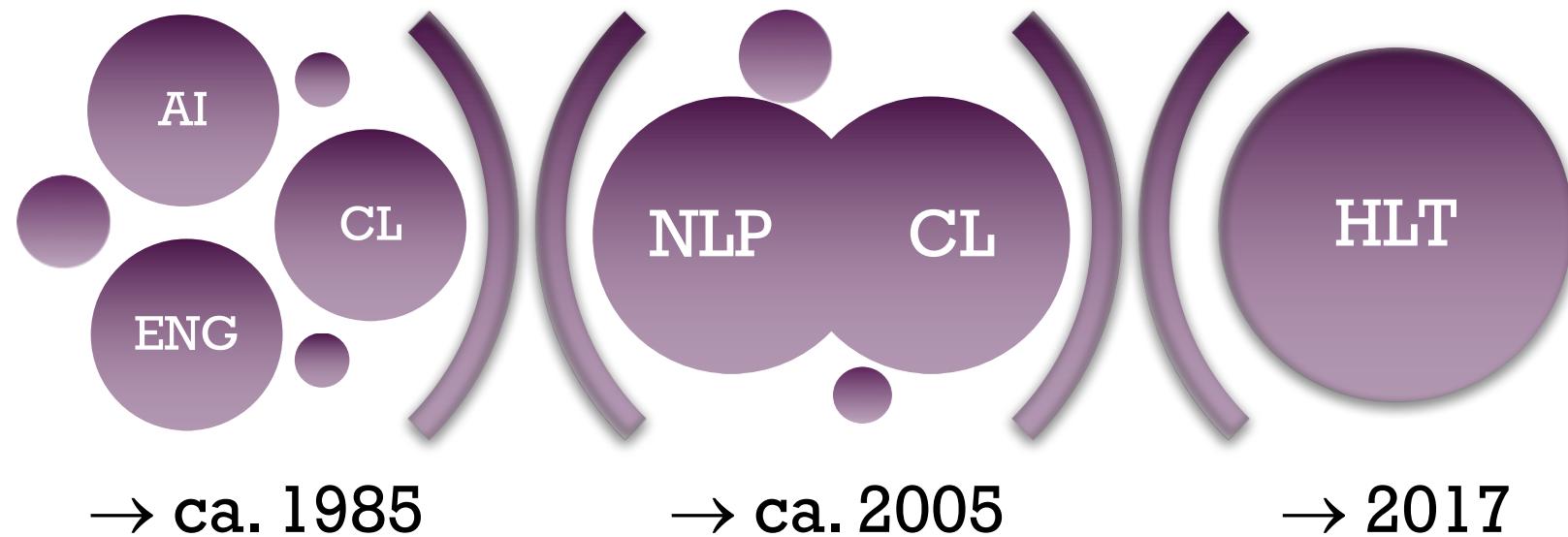
- To give an insight by way of some examples, rather than an exhaustive overview, into how the field of human language technologies has evolved. The focus today is on fundamental ideas with more specific elements covered in later lectures.
- Enrolment key on csmoodle.ucd.ie: 40020HLT19
(note 2018-19 version)

+ Some Literature

- Ittoo, A.; L.M. Nguyen & A. van de Bosch (2016): Text analytics in industry: Challenges, desiderata and trends, *Computers in Industry* (2016),
<http://www.sciencedirect.com/science/article/pii/S0166361515300646>
- Ledeneva, Y & G. Sidrov (2010): Recent Advances in Computational Linguistics, *Informatics* 34 (2010), 3-18.
- Kay, M (2006): A Life of Language, ACL Lifetime Achievement Award, *Computational Linguistics*, Volume 31, Number 4: 425-438.
- Cole, R. et al (1997): *Survey of the State of the Art in Human Language Technology*, Cambridge University Press & Giardini.

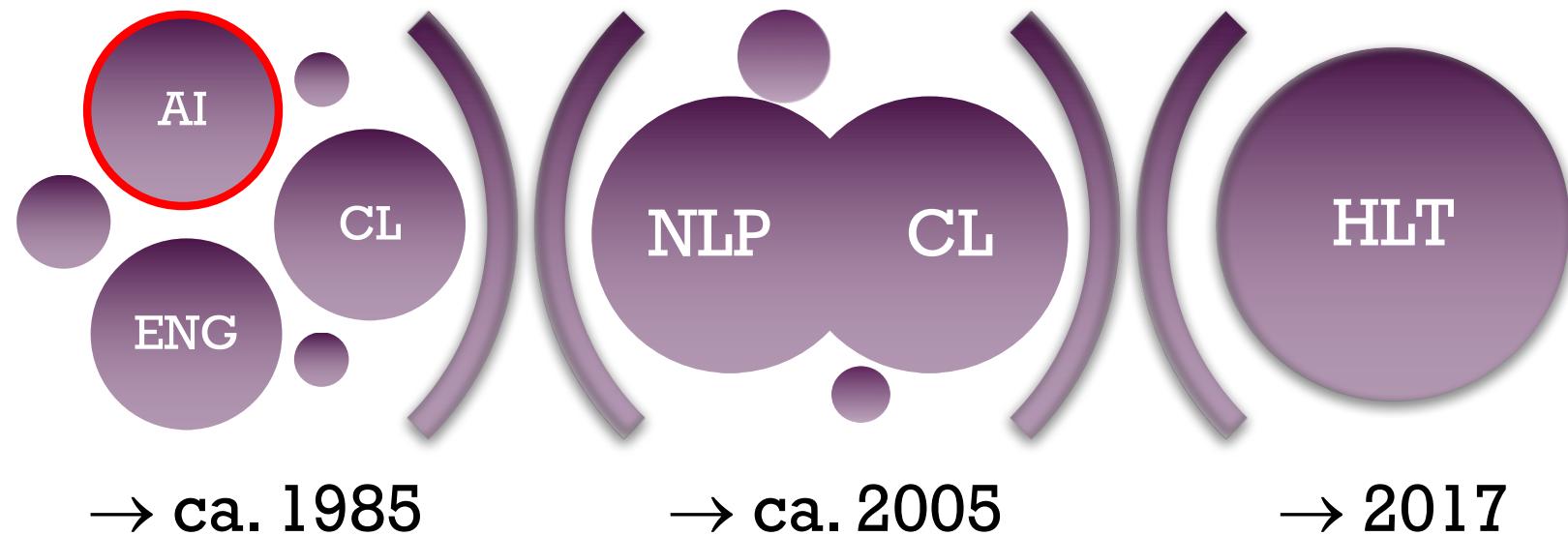
+ High-Level Timeline

HLT2



+ High-Level Timeline

HLT2



+ Early AI: ELIZA (1966)

HILT2

- ELIZA - models the behaviour of a psychiatrist (cf. Turing Test/Imitation Game)
- “The fundamental technical problems with which ELIZA is concerned are:
 - the identification of key words,
 - the discovery of minimal context,
 - the choice of appropriate transformations,
 - generation of responses in the absence of keywords, and
 - the provision of an ending capacity for ELIZA "scripts".

Joseph Weizenbaum (1966): ELIZA: A Computer Program For the Study of Natural Language Communication Between Man and Machine, *Communications of the ACM*, Volume 9, Number 1.

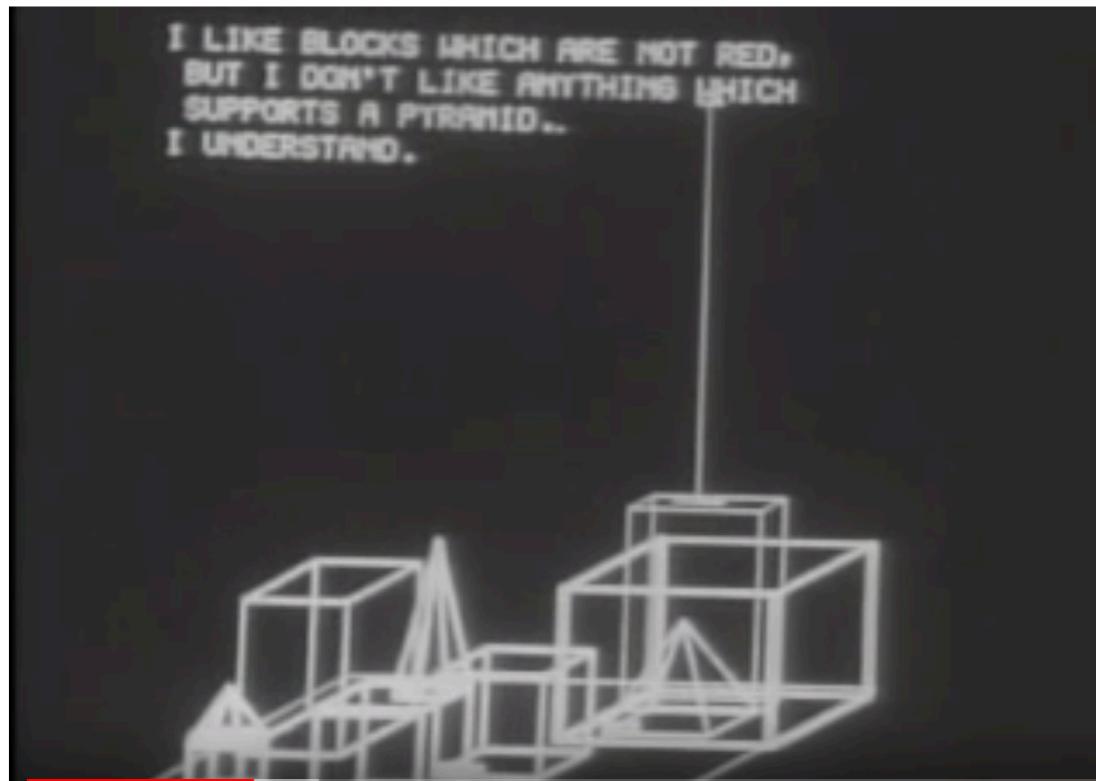
- See <http://www.masswerk.at/elizabot/> - written by Norbert Landsteiner
- Some more recent developments:
<http://www.bbc.com/news/technology-27762088>



+ Early AI: SHRDLU (1971/2)

HILT2

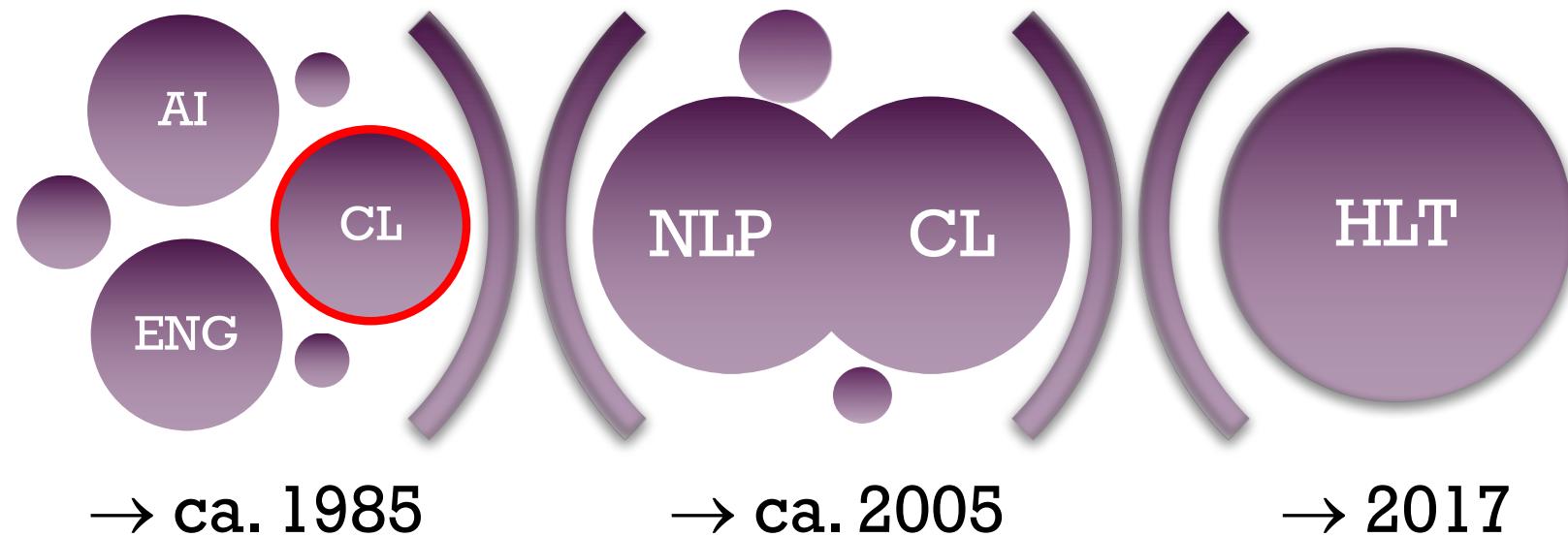
- SHRDLU – a language parser with natural language user interaction in a “blocks world” - Terry Winograd (1972): *Understanding Natural Language*, Academic Press.
- <http://hci.stanford.edu/winograd/shrdlu/>



- For a (silent) demonstration see:
<https://www.youtube.com/watch?v=bo4RvYJYOzI>

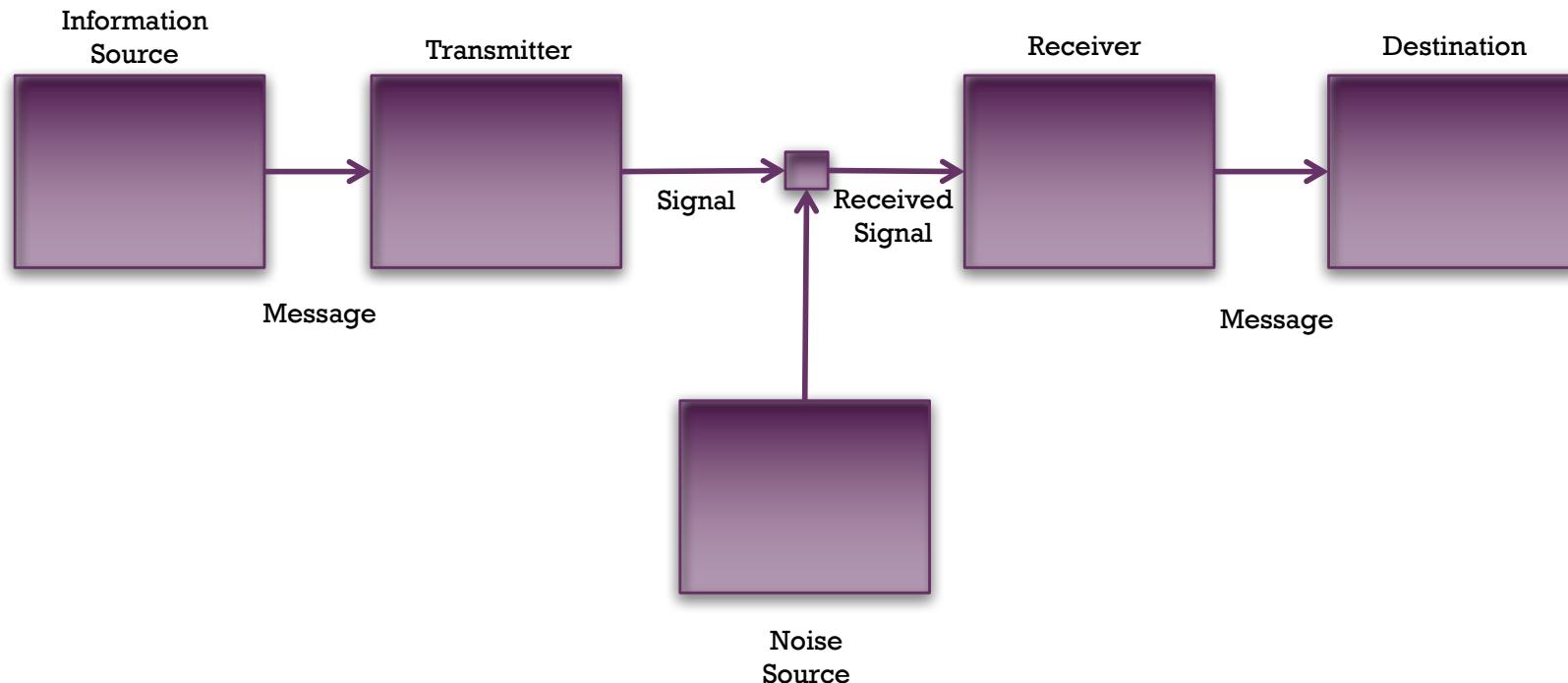
+ High-Level Timeline

HLT2



+ Early CL: Machine Translation

HILT2



Shannon's General Communication System (Shannon, 1948)

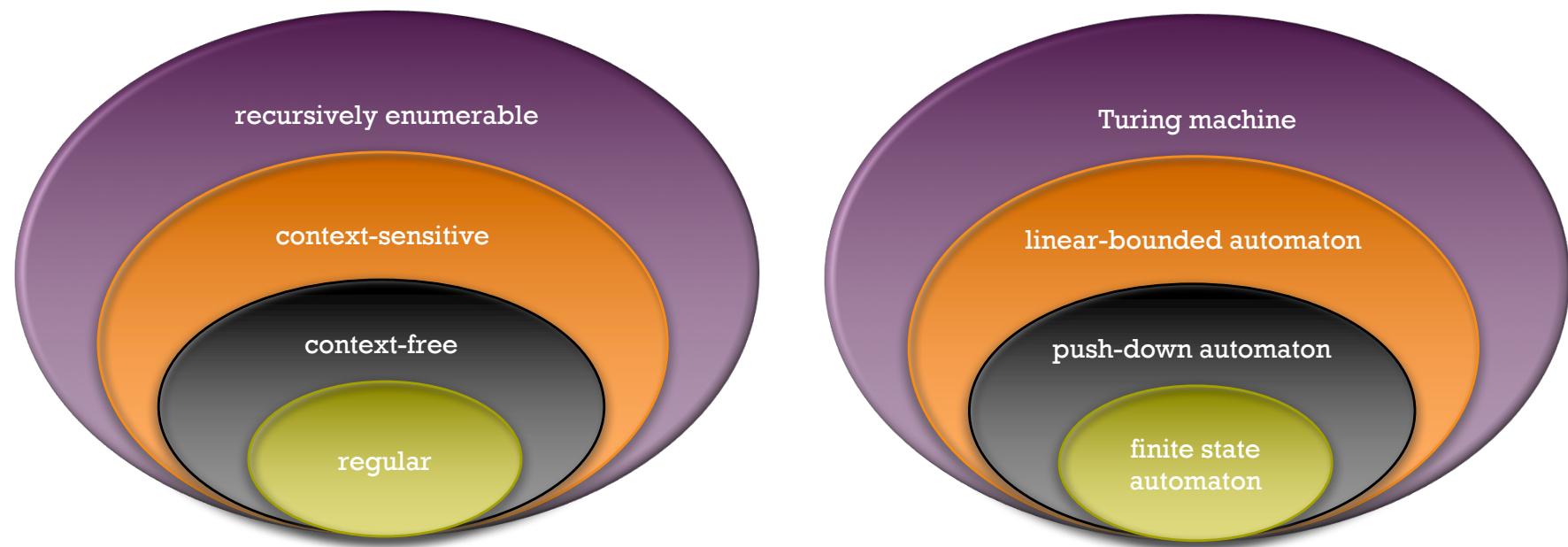
- Shannon, C. & W. Weaver (1949): *The Mathematical Theory of Communication*
- Weaver, W. (1949): *Translation*
- <http://www.mt-archive.info/Weaver-1949.pdf>



+ Early CL: “Formalising” Grammar

HLT2

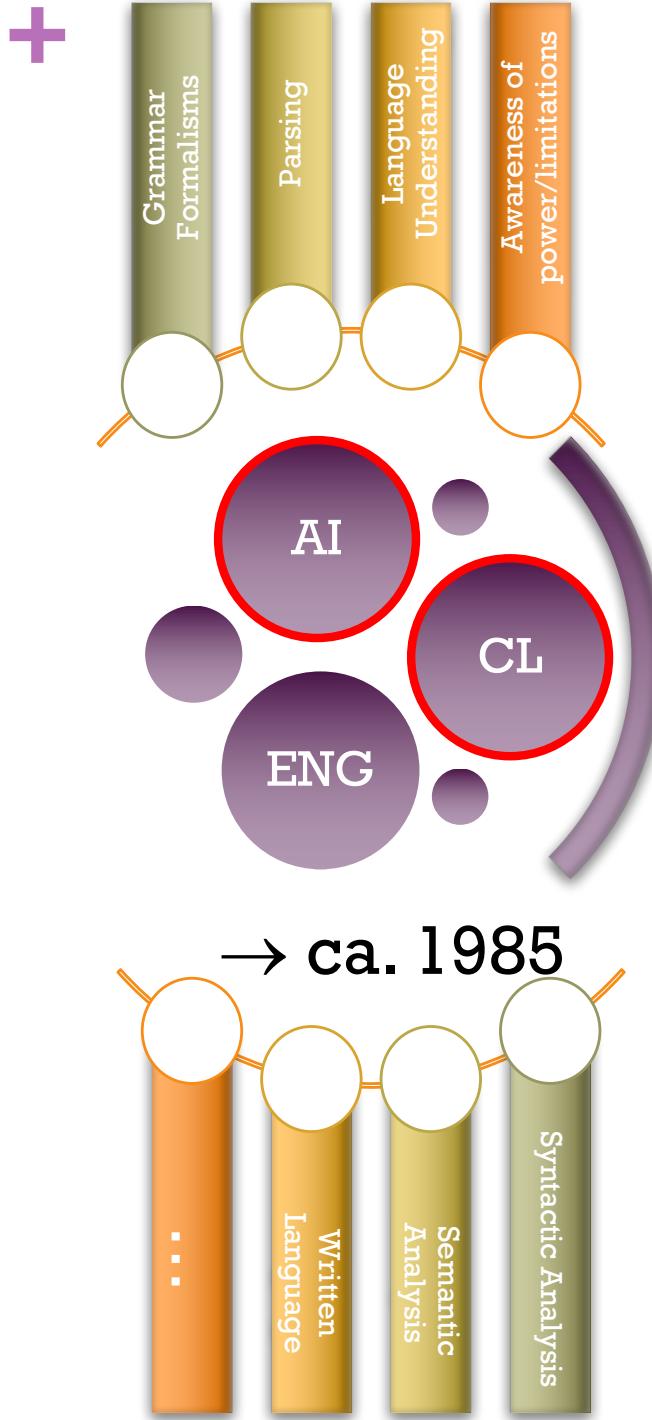
- Relating natural language syntax to formal language theory – mathematical models of language



Chomsky Hierarchy

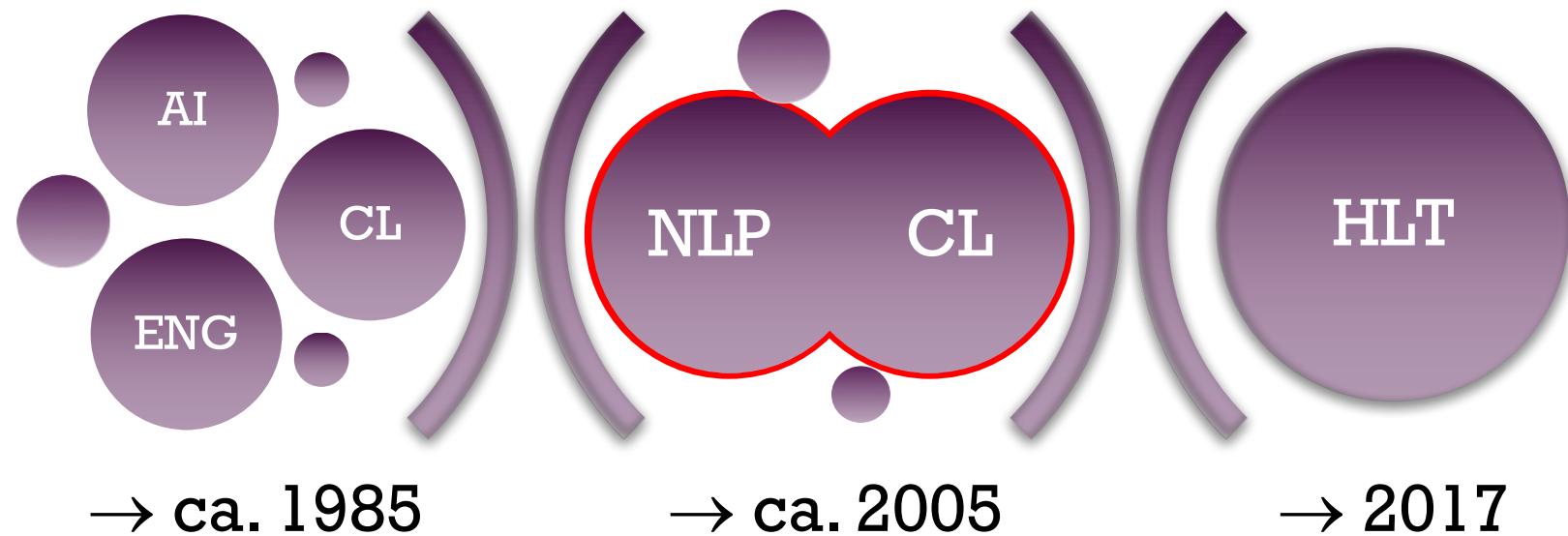


Some Examples



+ High-Level Timeline

HLT2



+ Cole et al. (1997) Review

HLT2

- The study of human language technology is a multidisciplinary enterprise, requiring expertise in areas of linguistics, psychology, engineering and computer science.
- Creating machines that will interact with people in a graceful and natural way using language requires a deep understanding of the acoustic and symbolic structure of language (**the domain of linguistics**), and the mechanisms and strategies that people use to communicate with each other (**the domain of psychology**).

Cole et al. (1997:xi)



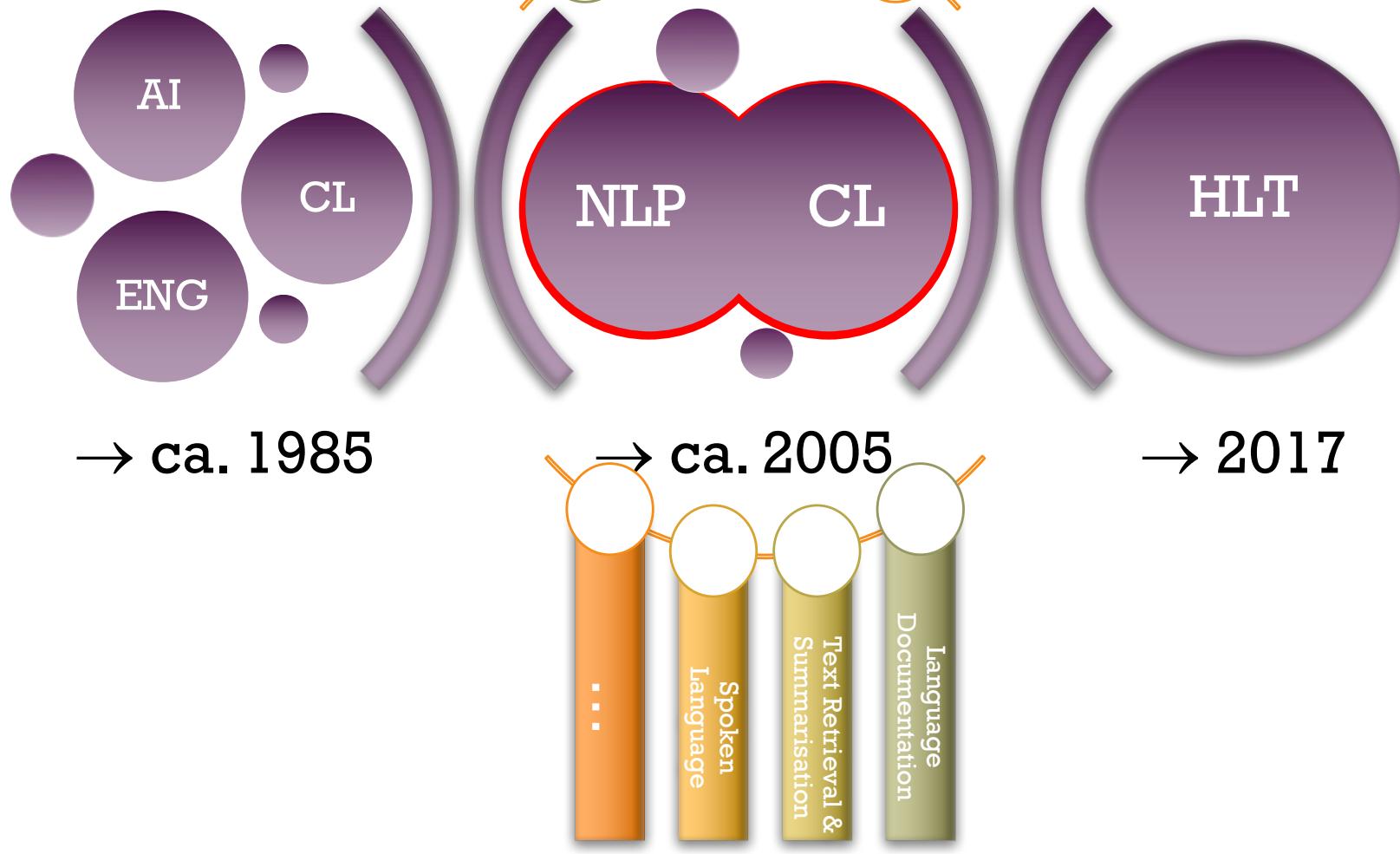
+ Cole et al. (1997) Review

HILT2

- Given the remarkable ability of people to converse under adverse conditions, such as noisy social gatherings or band-limited communication channels, advances in signal processing are essential to produce robust systems (**the domain of electrical engineering**).
- Advances in **computer science** are needed to create the architectures and platforms needed to represent and utilize all of this knowledge.
- Collaboration among researchers in each of these areas is needed to create multimodal and multimedia systems that combine speech, facial cues and gestures both to improve language understanding and to produce more natural and intelligible speech by animated characters.

+

Some Examples



+ (E)Merging “Dichotomies”...

HILT2

- Technological vs Linguistic/Cognitive
- Procedural vs Declarative
- Competence vs Performance
- Knowledge-Based vs Data-Driven
- Small (but difficult) Problems vs Broad Coverage
- Algorithms vs Resources
- Methods vs Areas of Research



+ Cole et al. (1997) Review

HILT2

Looking forward:

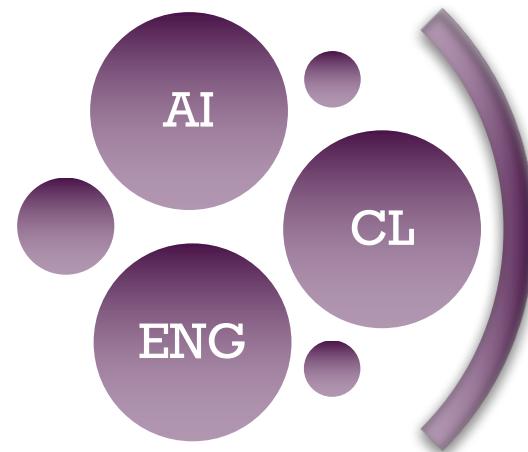
- Benefits that can be expected from deploying language technology are a more effective usability of systems (**enabling the user**) and enhanced capabilities for people (**empowering the user**). The economic and social impact will be in terms of efficiency and competitiveness for business, better educated citizens, and a more cohesive and sustainable society.

Cole et al. (1997:xvii)

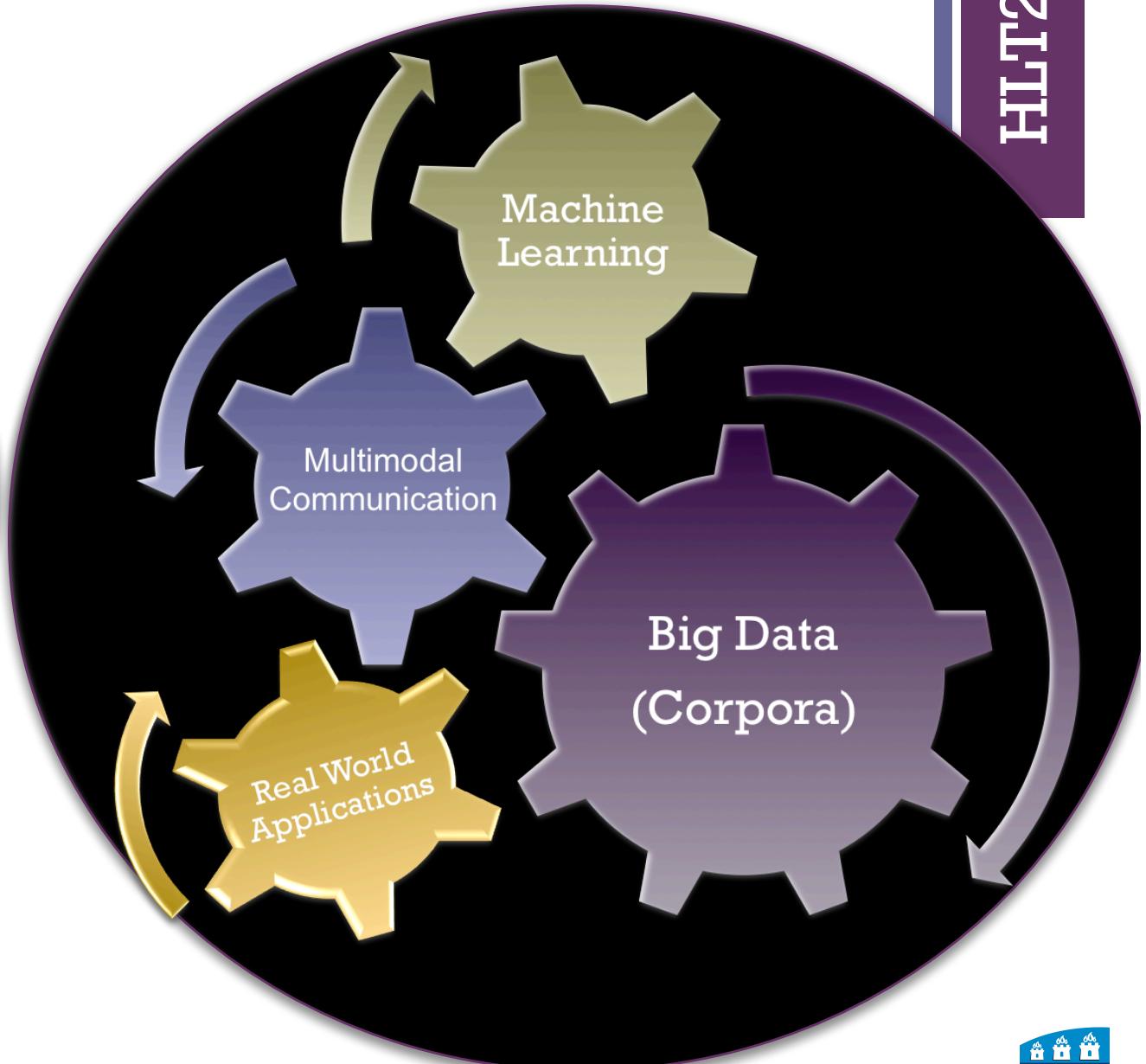


+

Focus



→ ca. 1985



+ Some Current HLT Applications

HLT2

- **Analysing and evaluating the task of automatic tweet generation: knowledge to business**
- “Text summarization techniques are applied to news articles, yielding ultra-concise summaries, which then serve as tweets. A potential application of these automatically generated tweets is in supporting the marketing and communication activities of companies. Therefore, it is important to assess the tweets’ quality, especially taking into account the fact that they are created without human intervention.
- The tweets’ quality is assessed in terms of 2 dimensions, viz:
 1. Interestingness: to what extent users would be interested in knowing more about the tweets’ contents.
 2. Informativeness: whether the tweets accurately reflected the news articles from which they were generated.”

E. Lloret, M. Palomar, **Analysing and evaluating the task of automatic tweet generation: knowledge to business**, Comput. Ind. (2015) according to Ittoo et al. (2016: 4)



+ Some Current HLT Applications

H2

- A methodology for traffic-related twitter messages interpretation
- “Automatically interpreting traffic-related Twitter messages in the Portuguese language. Interpretation in this case refers to the transformation of the tweets’ unstructured textual contents into a more structured formalism, namely as RDF-triples.
- The system is deployed as a prototype for monitoring the truck fleet of a gas transportation and fuel transportation company. Several benefits are reported as a result of the system’s application in these companies, such as cost reduction, more efficient fleet management, and improved customer satisfaction due to a more accurate prediction of delivery time.”

Resource
Description
Framework

F.C. Albuquerque, M.A. Casanova, H. Lopes, L.R. Redlich, J.A.F. de Macedo, M. Lemos, M.T.M. de Carvalho, C. Renso, A methodology for traffic-related twitter messages interpretation, Comput. Ind. (2015) according to Ittoo et al (2016: 4)



+ Some Current HLT Applications

HLT2

- **Integrating a semantic-based retrieval agent into case-based reasoning systems: a case study of an online bookstore**
- “employs Case-Based-Reasoning (CBR) to enhance the search experience of end users in business-to-consumer websites. Specifically, it accepts queries as inputs, which could be formulated either as natural language questions or as keywords. It then consults a case base consisting of target problems (e.g. questions) and their corresponding solutions (answers). The best matching solution in the case is then returned as answer to the input query.”

J.W. Chang, M.C. Lee, T.I. Wang, **Integrating a semantic-based retrieval agent into case-based reasoning systems: a case study of an online bookstore**, Comput. Ind. (2015) according to Ittoo et al (2016: 6)



+ Some Current HLT Applications

HLT2

- A distributional approach to open questions in market research
- “system called The Klugator Engine (TKE), for analyzing responses to open-ended survey questions. The responses are expressed in freely formed natural language texts in English and German. The proposed system was developed in a collaborative effort between academic and industrial partners. The system is already in commercial use as the core engine of the Rogator Text Clustering Solution (RogTCS), distributed by German market research company Rogator AG.”

P. Greiner, S. Evert, F. Baigger, B. Lang, A distributional approach to open questions in market research, Comput. Ind. (2015) according to Ittoo et al (2016: 6)



+ Some Current HLT Applications

HLT2

- IBM Watson – a technology platform that uses natural language processing and machine learning to reveal insights from large amounts of unstructured data
- <http://www.ibm.com/smarterplanet/us/en/ibmwatson/what-is-watson.html>
- Winning at Jeopardy:
https://www.youtube.com/watch?v=WFR3lOm_xhE



+ NLP Industry...

HLT2

“The Natural Language Processing (NLP) market size is estimated to grow from USD 7.63 Billion in 2016 to USD 16.07 Billion by 2021, at a Compound Annual Growth Rate (CAGR) of 16.1%.”

marketsandmarkets.com, July 2016 (Report Code: TC , 3492)

A huge
part of AI



+ NLP Industry...

HLT2

GENERAL

what industries are next
to be disrupted by NLP
and Text Analysis?

September 16, 2016 - General

f Facebook

Twitter

in LinkedIn



<http://blog.aylien.com/nlp-text-analysis-insurance-legal-customer-service/>

+ NLP Industry...

How AI and NLP will affect the media industry

One of the most visible of these technologies has been the use of artificial intelligence (AI) as well as Natural Language Processing (NLP), in process of collecting data and analyzing it.

ETtech | Updated: October 08, 2017, 21:13 IST

[Share 21](#) [G+ Share](#) [in Share](#) 58 [Tweet](#)

By Sabir Chowdhury, Director, Global Technology Services at Time Inc. India

Publishing industry across the world is going through challenges. The dynamic nature of technology trends demands its continuous evolution from publishing to a digital media company. Progress has been made in terms of both content platforms i.e. the move from purely print to a variety of audio-visual avenues (such as television and online news portals, among others), as well as in terms of technology used to gather and publish information. One of the most visible of these technologies has been the use of artificial intelligence (AI) as well as Natural Language Processing (NLP), in process of collecting data and analyzing it.



<https://tech.economictimes.indiatimes.com/news/technology/how-ai-and-nlp-will-affect-the-media-industry/60966371>

+ NLP Industry...

HILT2

sort by: Relevance

11d

13d

16d

11d

Linguist, Text Classification
Google
Dublin, IE
Analytical Linguists work in many different areas and arrive with a wide variety of skills—your specialization might involve natural language processing and understanding, ...
careers.google.com

Watson Cognitive Engineer
IBM
Dublin, IE
We are looking for people who has skills about machine learning, NLP, information retrieval, named entity recognition, word-sense disambiguation, language modeling, parsing, ...
sjobs.brassring.com

Data Analyst with English (UK Market)
Globe Tech Inc.
Cork, IE
Keywords: Jobs in London, Jobs in Cork, Graduate jobs, Bilingual jobs, Jobs with languages, Tech support jobs, IT jobs, Temporary jobs, Contract jobs, Bilingual jobs, Jobs with ...
careers.globetech.icims.com

Senior Data Scientist
Ancestry
Dublin, IE
We're looking for a Senior Data Scientist to join and help grow this exciting team. 4+ years working experience in machine learning / statistical modeling.
jobs.smartrecruiters.com

How AI is changing publishing

One of the most visible ways AI is changing publishing is through Natural Language Processing (NLP). This technology allows computers to understand, interpret, and generate human language, making it easier for publishers to analyze and process large amounts of text.

By Sabir Chowdhury, Director, Publishing industry across the company. Progress has been made purely print to a variety of audio-portals, among others), as well as intelligence (AI) as well as Natural Language Processing jobs, careers, employment in Ireland

September 2018

Facebook



<https://ie.linkedin.com/jobs/natural-language-processing-jobs>

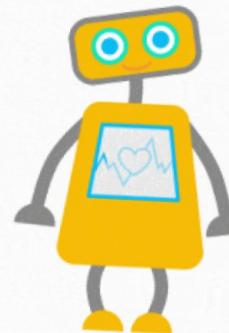
+ NLP Industry...

HIT2

To Your Health

'The Woebot will see you now' – the rise of chatbot therapy

By Amy Ellis Nutt December 3, 2017



Woebot is always available and will never judge. (

My therapist wanted to explain a few things during our first online session:

"I'm going to check in with you at random times. If you can't respond straight away, don't sweat it. Just come back to me when you're ready. I'll check in daily."

Alison Darcy CEO
Woebot Labs and
formally UCD
Psychology



https://www.washingtonpost.com/news/to-your-health/wp/2017/12/03/the-woebot-will-see-you-now-the-rise-of-chatbot-therapy/?utm_term=.1c5b400d32db

+ Speech Industry...

HILT2

How voice technology is transforming computing

Like casting a magic spell, it lets people control the world through words alone

Print edition | Leaders >
Jan 7th 2017

Normal

Twitter Facebook LinkedIn Email Print Comment

Computers must be able to understand context in order to maintain a coherent conversation about something, rather than just responding to simple, one-off voice commands, as they mostly do today ("Hey, Siri, set a timer for ten minutes").

<https://www.economist.com/news/leaders/21713836-casting-magic-spell-it-lets-people-control-world-through-words-alone-how-voice>



+ Speech Industry...

HILT2

BIG DATA

The future of voice synthesis after Google WaveNet debut

Dr. Matthew Aylett Tue 13 Sep 2016 1.03pm



G+ 3 Twitter 26 Y 1 f 9 in 3 Reddit 37 < 78 SHARES

Dr Matthew Aylett is the Chief Scientific Officer and co-founder at CereProc, which develops voice synthesis systems primarily for the healthcare field. In the light of last week's announcements from Baidu/NVidia and from Google, he takes a look at the challenges and motivations involved in reproducing the human voice artificially...

When Google first took Siri out of the box they didn't say "Nice recognition, interesting application, cool array microphones" – they said "Where are our Ads?" Google's business model almost entirely depends on selling Ads. Siri replaces the browser and returns control to Apple.



<https://thestack.com/big-data/2016/09/13/the-future-of-voice-synthesis-after-google-wavenet-debut/>



+ Speech Industry...

HILT2

The image shows a collage of various news snippets and social media posts related to AI speech synthesis. At the top left, a large headline from Forbes reads "How voice tech... This Startup's Artificial Voice Sounds Almost Indistinguishable From A Human's". Below it, another snippet from Forbes by Parmy Olson discusses an Irish startup's breakthrough in text-to-speech synthesis, mentioning DeepMind and Facebook. To the right, a red speech bubble contains the text "Peter Cahill, CEO of VOYSIS and formerly UCD Computer Science". At the bottom, a snippet from a news site quotes Peter Cahill as saying "you no longer need a multi-billion dollar R&D budget or hundreds of engineers to produce an artificial voice that's as good as Google's". The collage also includes small images of people and social media sharing icons.

<https://www.forbes.com/sites/parmyolson/2017/11/03>this-startups-artificial-voice-sounds-almost-indistinguishable-from-a-humans/#12fed267388c>



+ Speech Industry...

HILT2

The image shows a BBC News article titled "Adobe Voco 'Photoshop-for-voice' causes concern". The article is dated 7 November 2016 and is categorized under Technology. It features a photograph of a stage where two people are sitting at a table, with a large screen behind them displaying a waveform from a speech recognition or synthesis application. The text of the article discusses the ethical and security concerns raised by the new application.

An Irish startup has created a breakthrough that improves speech recognition by Google.

The results have caused many concerns from Amazon and several hundred Google researchers.

Voice recognition is

A new application that promises to be the "Photoshop of speech" is raising ethical and security concerns.

<http://www.bbc.com/news/technology-37899902>





Module Outline

HLT2

- HLT 1: Introduction
- HLT 2: Evolution of Language Technologies
- HLT 3: Workshop 1 (Monday & Wednesday)
- HLT 4: Word-forms, wordforms and word forms (**Morphology**)
- HLT 5: Colourless green ideas sleep furiously (**Syntax, Semantics**)
- HLT 6: Grammars, Rules & Parsers (**Syntactic Parsing**)
- HLT 7: Workshop 2 (Wednesday & Monday)
- HLT 8: Boundaries, Tokens & Corpora (**Text Processing & Analytics**)
- HLT 9: Same but Different (**Similarity & Sentiment Analysis**)
- HLT 10: Workshop 3 (Wednesday & Monday)



+ Module Outline

- HLT 11: Sounds, Symbols & Features (**Phonetics & Phonology**)
- Break
- HLT 12: Cats, Flys & Tomatos? (**Computational Morphology**)
- HLT 13: Twas brillig and the slithy toves (**Computational Phonology**)
- HLT 14: Language Modelling
- HLT 15: Workshop 4 (Wednesday & Monday)
- HLT 16 Speech Recognition
- HLT 17: Speech Synthesis
- HLT 18: Applications



HLT2

+ Jupyter Notebook & NLTK

NLTK2

- For Monday, please install Jupyter Notebook and NLTK by following the instructions on csmoodle
- If you have any difficulties with the installation, please contact me by email
- Please bring a laptop on Monday

