

We now ask the question, 'What will happen when a machine takes the part of A in this game?'

...

These questions replace our original, 'Can machines think?'

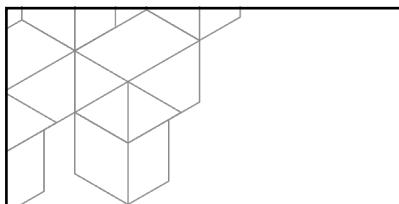
Alan Turing

COMP6713 | 2025 T1 | UNSW

All images from Wikimedia commons, unless specified.



1



COMP6713 Natural Language Processing (NLP)



Convener

Dr. Aditya Joshi

aditya.joshi@unsw.edu.au



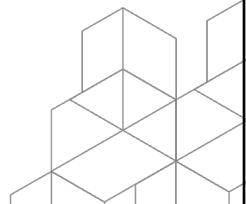
Week 1

Introduction



Schedule

2025 Term 1



2

Course Convener

Name: Aditya Joshi (he/him)

Office: K17-217B

Email: aditya.joshi@unsw.edu.au

Questions?

Technical: WebCMS Forum

Technical/personal: Email

Hobbies: Bushwalking, cooking, standup comedy [[Sydney Fringe '24](#)]

Research: Foundational and Applied Natural language processing

(<https://unswnlp.github.io/>)



Youtube screenshot



3

Course Convener

Name: Aditya Joshi (he/him)

Office: K17-217B

Email: aditya.joshi@unsw.edu.au

Questions?

Technical: WebCMS Forum

Technical/personal: Email

Hobbies: Bushwalking, cooking, standup comedy [[Sydney Fringe '24](#)]

Research: Foundational and Applied Natural language processing
(<https://unswnlp.github.io/>)

NLP and me:

- Masters Project (2010): Sentiment analysis of tweets

- PhD Thesis (2018): Investigations in Computational Sarcasm

Experience

- CSIRO's Data61: NLP for epidemic intelligence

- SEEK: NLP for candidate-job recommendation

- Notiv (startup acquired by Dubber): NLP for meeting summarisation



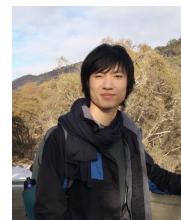
4

2

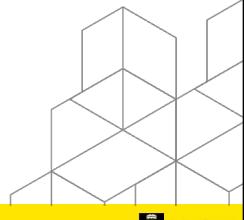
Course Team



Lihua Wang (Lily) Martin Eftimoski

Rashini
Liyanarachchi
(Rashini)Duc Anh
Nguyen (Duke)

Zechen Li (Ryan)



Photos from LinkedIn

5

 UNSW SYDNEY | Australia's Global University

Module 1
Introduction

Course Overview

- Objectives
- Content
- Assessments

Layers of NLP

Introduction to NLP tasks
Ambiguity, probability and data.

Black-box NLP

NLP Libraries: Spacy, NLTK,
HuggingFace pipelines
Key considerations for NLP.
...

6

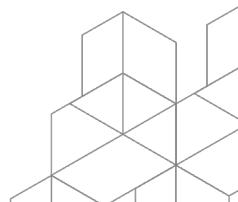


7

COMP6713: Natural Language Processing

- Introductory course; second offering
- Philosophy:
 - What: Language phenomena
 - How: Computational models
 - Why: Historical background

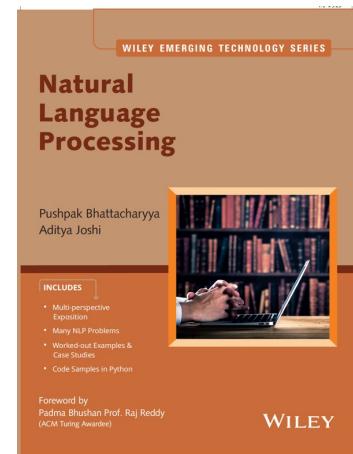
Is this course an advanced disciplinary elective (or similar term) for you? Please contact the Student Nucleus Hub.



8

Recommended Resources

- 1) Daniel Jurafsky and James H. Martin. 2025. Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition with Language Models, 3rd edition. Online manuscript released January 12, 2025. <https://web.stanford.edu/~jurafsky/slp3>.
- 2) Pushpak Bhattacharyya and Aditya Joshi. 2023. Natural Language Processing. Kindle Edition released December 2023. Wiley. <https://www.amazon.com.au/Natural-Language-Processing-Pushpak-Bhattacharyya-ebook/dp/B0CR64RX4T> (Copies available in the Library)



9

Pre-requisites

Postgraduate: COMP9020 and (COMP9814 or COMP9444)*

Undergraduate: MATH1081 and (COMP3411 or COMP9444)*

Will I enjoy the course? (Not a complete list. Not necessary conditions)

- I am a language nerd
- I enjoy coding
- I like learning about AI models
- I often connect modern ideas with those from the past

*Exceptions made.



10

Learning Objectives

- CL01: Describe NLP problems such as POS tagging, sentiment analysis, information extraction and machine translation along with their challenges in terms of ambiguity resolution.
- CL02: Explain typical NLP approaches based on statistical and neural approaches.
- CL03: Use NLP libraries (e.g. NLTK, scikit-learn, Transformers) to implement the training of models for NLP problems and use them for inference.
- CL04: Design an NLP solution by selecting the NLP problem formulation, approach and evaluation strategy, by analysing the requirements of a specific application.



11

Why are we even studying pre-deep learning NLP?

- Large language models (LLMs) are far from perfect
- Pre-deep learning NLP approaches are a ‘small’ component of the course
- They are still used in several businesses.. in Australia
- Concepts from pre-deep learning NLP models help us appreciate the ingenuity of deep learning-based NLP
 - This goes to the last slide:
 - Linguistic examples help us understand the ‘what’ of NLP
 - Neural models help us understand the ‘how’ of NLP
 - Three-generational view helps us understand the ‘why’ of NLP

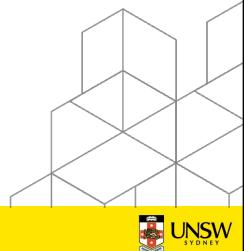


12

How is the 2025T1 offering different from 2024T1?

Student feedback:

- Overall: Strongly Positive
- Appreciated the engagement and humour in the class.
- Agreed that it was a valuable course.

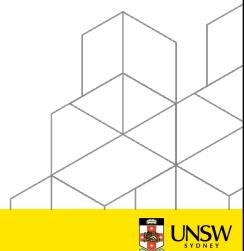


13

How is the 2025T1 offering different from 2024T1?

Student feedback:

- Overall: Strongly Positive
- Appreciated the engagement and humour in the class.
- Agreed that it was a valuable course.
- Thought the assignment was not handled sufficiently well.
- Thought more engagement with the lecturer would have helped.
- Drag-and-drop questions in Moodle don't work well on all devices.



14

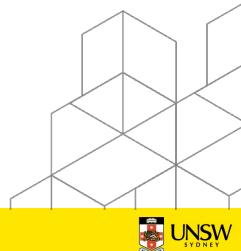
How is the 2025 T1 offering different from 2024 T1?

Student feedback:

- Overall: Strongly Positive
- Appreciated the engagement and humour in the class.
- Agreed that it was a valuable course.

- Thought the assignment was not handled sufficiently well.
 - Revised assignment with clear evaluation criteria.
- Thought more engagement with the lecturer would have helped.
 - Option of online consultation in addition to in-person consultation.
- Drag-and-drop questions in Moodle don't work well on all devices.
 - Cool, None.

- Also, new content in EVERY module to keep up with the fast-changing NLP landscape.
- New format for group project (industry project option) and final exam (Inspera – TO BE CONFIRMED).
- etc.



15

Content Plan (Subject to Minor Revision)

Week#	Module	Topics
1	Introduction	NLP tasks, ambiguity resolution and generations of NLP, Ethical considerations, Black-box NLP libraries; Datasets + API calls for LLMs
2	Representation learning	Grammar, Probabilistic language models, Word Vectors + Review of Sequential models
3	Attention & Transformer	Attention, Fine-tuning, Prompt tuning
4	Language models	Encoder models, decoder models, LoRA, LangChain +....
5	Sentiment Analysis	Lexicons, Statistical classifiers, LSTM stacks, BERT-based models+ Evaluation, Benchmarks
7	POS tagging and NER	HMM-based POS tagging, POS tagging a seq2seq task, CRF-based NER
8	Machine translation	Rule-based, statistical, Neural, Decoding+
9	Summarisation & Question-Answering	Extractive and abstractive summarization +
10	Other NLP tasks	Shared tasks and benchmarks +Detailed benchmarks; NLP applications

+ Guest lectures; “How to read an NLP paper” and so on.



16

Course Schedule

Lectures:

Thursday: 14:00-16:00. O'Shane 104

Friday: 11:00-13:00. Griff M18

Public Holiday (Good Friday): Friday 18 April 2025

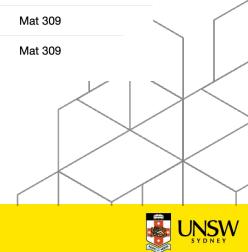
Have you been added to WebCMS and Moodle?

Have you received your tutorial allotments?

<https://webcms3.cse.unsw.edu.au/COMP6713/25T1>

Tutorials

Name	Staff	Day	Start Time	End Time	Weeks	Room
H17A		Thu	17:00	18:00	1-5,7-10	Myers 540
		Thu	17:00	18:00	1-5,7-10	Myers 540
H17B		Thu	17:00	18:00	1-5,7-10	Myers 440
		Thu	17:00	18:00	1-5,7-10	Myers 440
F10B		Fri	10:00	11:00	1-5,7-8	Quad G026
		Fri	10:00	11:00	1-5,7-8	Quad G026
F13A		Fri	13:00	14:00	1-5,7-8	Law 276
		Fri	13:00	14:00	1-5,7-8	Law 276
F13B		Fri	13:00	14:00	1-5,7-8	Law 203
		Fri	13:00	14:00	1-5,7-8	Law 203
F14A		Fri	14:00	15:00	1-5,7-8	Mat 309
		Fri	14:00	15:00	1-5,7-8	Mat 309



17

Assessments (1/2)

Weekly Quizzes (20%):

- Moodle-based, closed-book quizzes
- Timeline:
 - Weekly: 2-5, 7-9
 - One attempt; no time limit
- Evaluation: Auto-marked
- Quiz will be open from Friday 5:25pm to Wednesday 5:25pm | SAMPLE quiz this week

Assignment (10%):

- Individual programming assignment
- Timeline:
 - Released in Week 3
 - Due in the end of week 5
 - Based on content from Weeks 1 to 3
- Evaluation:
 - Test scripts
 - Manual evaluation (style; correctness, etc.)



18

Assessments (2 / 2)

Group Project (25%):

- Groups of 3 – 5
- Timeline:
 - Week 4: Register your group and topic in consultation with the team
 - Project due on Monday of Week 10
- Deliverables:
 - Code, Report, Presentation/video

Final Exam (45%):

- 2 hours; during exam period
- Multiple-choice questions, short-answers and code analysis.

To pass this course, students must score more than 40% on the final exam.



Assessments (2 / 2)

Group Project (25%):

- Groups of 3 – 5
- Timeline:
 - Week 4: Register your group and topic in consultation with the team
 - Project due on Monday of Week 10
- Deliverables:
 - Code, Report, Presentation/video + In-person/Online Q&A

Final Exam (45%):

- 2 hours; during exam period
- Multiple-choice questions, short-answers and code analysis.

To pass this course, students must score more than 40% on the final exam.



Exciting news this year!

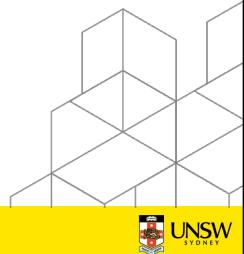
You may have the option of working
on an industry-led group project.



Integrity

You acknowledge that you understand the [UNSW Academic Integrity policy](#).

Also, a reminder about [Academic Integrity in Programming Courses](#)



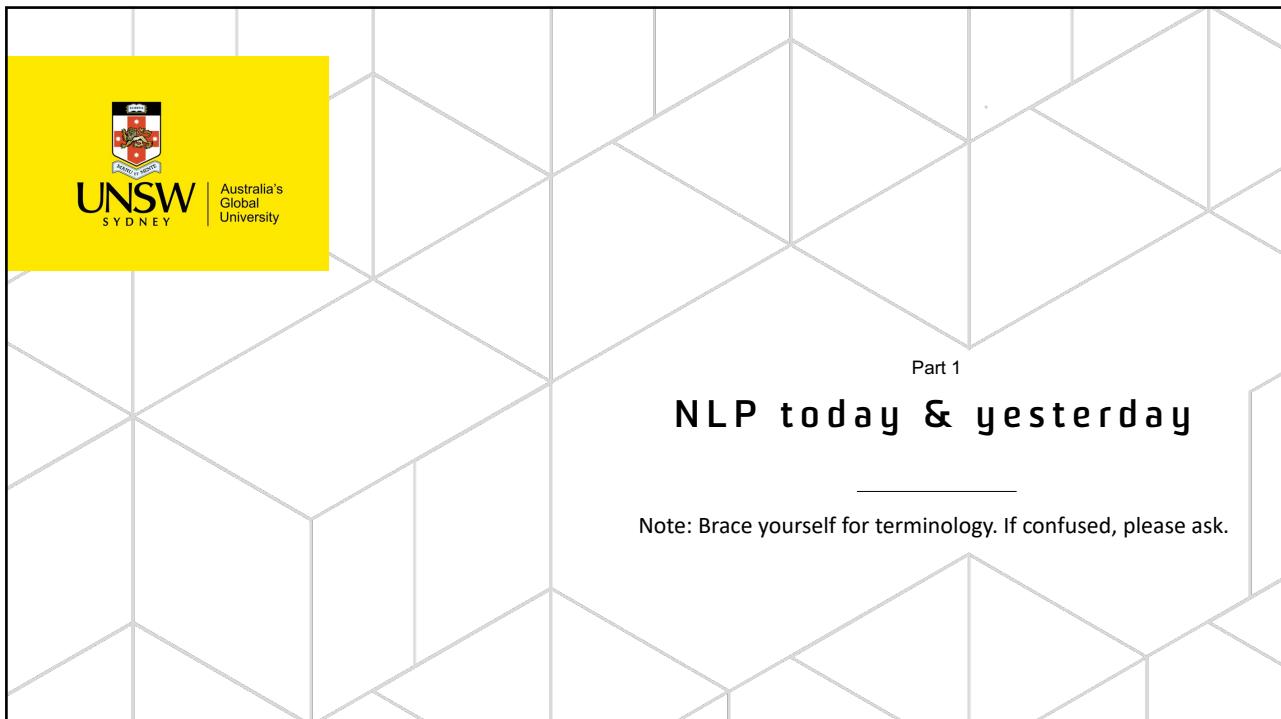
21

Course Resources

- Lecture Slides & Recordings
- Tutorial Problem Sets (Not evaluated)
- Suggested textbooks:
 - *Natural Language Processing*, Pushpak Bhattacharyya and Aditya Joshi, Wiley, 2023.
 - *Speech and Language Processing*, Daniel Jurafsky and James Martin, Prentice Hall, 2023.



22



The slide features the UNSW Sydney logo in the top left corner, consisting of a crest and the text "UNSW SYDNEY Australia's Global University". The background is a light gray with a complex pattern of intersecting triangles and lines forming a grid-like structure.

NLP today & yesterday

Note: Brace yourself for terminology. If confused, please ask.

23

What is natural language processing?

Natural language processing (NLP) is the branch of artificial intelligence that deals with computational processing of text.

Also known as 'computational linguistics'..

..and has a significant overlap with 'human language technology'.... or 'text analytics'.

- Natural languages: Language that humans use
- Artificial languages: Programming languages



24

NLP Today

The quick brwn fox jumps over the lazy dog

Change to: brown

Bren
Brown
brown
brown
brown

deepseek

grammarly

"They are buried underground and they do pose a small risk to firefighters."

He said the Australian Defence Force had been working with the RFS to advise them on the exact ex...

Evans Head Heritage Copy

Copy Link to Highlight

Search Google for "exclusion"

Print

Translate Selection to English

Richard Gates said despite all

exploded.

If it was anything like Evans Head, I can tell

1>Password

Inspect

Speech

Services

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

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

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

https://en.wikipedia.org/wiki/Spell_checker

https://www.business-standard.com/world-news/deepseek-rl-chinese-ai-research-breakthrough-challenging-onenet-explained-125012700327_1.html

25

NLP finds several applications

Headache

Fever

Cough

Diarrhoea

Vomiting

Sore throat

Pain

Runny nose

Post course

Signal indicator

Recent case

Report within a month

Report within 1 day

VIC

Post course on map

200 200

Figure 1: An overview of tasks in LegalAI.

LEGITIMATE	FAKE
<p>Nintendo Switch game console to launch in March for \$299. The Nintendo Switch video game console will sell for about \$299 in North America with 3 million units expected to be sold during its global rollout in the U.S. and Europe. The Japanese company promises the device will be packed with fun features of all its past machines and more. Nintendo is promising a more immersive, interactive experience with the Switch, including online playing and using the remote controller in games that don't require players to be constantly staring at a display.</p>	<p>New Nintendo Switch game console to launch in March for \$299. Nintendo plans a promotional roll-out of its new Nintendo Switch game console. For the first time, the console will roll out for an introductory price of \$99. Nintendo promises to pack the new console with fun features not present in past machines. The new console contains new features such as motion detectors and immersive and interactive gaming. The new introductory price will be available for two months to show the public the new advances in gaming.</p>

Table 2: Sample legitimate and crowdsourced fake news in the Technology domain

Brian Liang, Aditya Joshi*, Ross Spears, Stephen Man, Cecile Paris and C Raine MacIntyre, "Watch The Flu: A Twitter Monitoring Tool for Epidemic Intelligence of Influenza in Australia", AAAI 2020.

Hanlin Zhang, Zhenyu Guo, Junming Chen, Tianqi Tang, Zhiyong Wang, and Jiaxin Liu, "A Deep Learning Model for News Bots NLP Recognition System: A Summary of Legal Artificial Intelligence", In ACL 2020.

Verónica Pérez-Rosat, Bennett Kleinberg, Alexandra Lefteris, and Radu Mihalcea, 2018. Automatic Detection of Fake News. In COLING.

26

Early NLP

NLP has long explored automation of natural language:

1964: ELIZA: Chatbot: Counsellor

1968: SYSTRAN: Translation system (Russian to English translation)

2006: IBM Watson: Question-answering system

Eliza, the Rogerian Therapist

ELIZA is a computer program that emulates a Rogerian psychotherapist. Just type your questions and concerns and hit return. Eliza will answer you.

```
Talk to Eliza
> Hello, I am Eliza.
* How do you know Eliza?
> Why are you interested in whether or not I am Eliza?
* Whoa, I thought you were an emulation of a
psychotherapist.
> We were discussing you, not me.

Input: 
```

When the original ELIZA first appeared in the 60's, some people actually mistook her for human. The illusion of intelligence works best, however, if you limit your conversation to talking about yourself and your life.

This javascript version of ELIZA was originally written by [Michel Wallace](#) and significantly enhanced by [George Dunlop](#).

© 1999 manifestation.com [PMB 417 - 10800 Alpharetta Hwy, Suite 208 - Roswell, GA 30076]
webmaster: sabren@manifestation.com

<https://psych.fullerton.edu/mbirnbaum/psych101/eliza.htm>



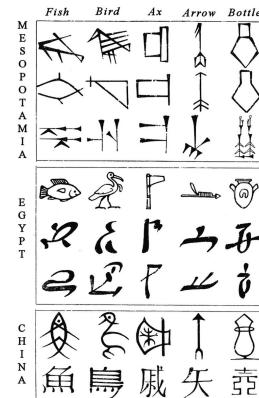
Natural Language

- Languages that humans speak (Examples: Mandarin Chinese, Spanish, English, Arabic, Hindi, Bengali)
- Language serves multiple functions: share knowledge, collaborate, .. Gossip
- Language communicates ideas
- Some 'languages' are specialised:
 - Language of a doctor
 - Language of an electrician
 - Language of a teacher teaching NLP
 - Language of a teacher teaching law
- Languages are different:
 - Native versus non-native speakers
 - Dialects versus so-called 'standard' language
 - Register (politeness versus argument)
- Language is an important factor for the success of the human race



Words

- Fundamental unit of a language
- Sounds produced to convey an idea
- Sounds were later written down
- Written scripts allowed ideas to be passed down to generations
 - Hieroglyphics
- A word conveys an idea
 - 'apple' (English) ->
 - 'say-b' (Hindi) ->
 - 'úll' (Irish) ->
- 'apple' (English)



https://en.wikipedia.org/wiki/Egyptian_hieroglyphs

Images from Wikimedia commons



29

.. The curious case of tea/cha

Derivatives of *te* [edit]

Language	Name	Language	Name	Language	Name	Language	Name	Language	Name
Afrikaans	<i>tee</i>	Armenian	պէ [pɛ]	Basque	<i>tea</i>	Belarusian	рапідэра [rapidiera]	Berber	ت, atay
Catalan	<i>te</i>	Kashubian	(n)arbatta ₍₁₎	Czech	<i>te</i> or <i>thé</i> ₍₂₎	Danish	<i>te</i>	Dutch	<i>thee</i>
English	<i>tee</i>	Esperanto	<i>teo</i>	Estonian	<i>tee</i>	Faroese	<i>te</i>	Finish	<i>tee</i>
French	<i>thé</i>	West Frisian	<i>te</i>	Galician	<i>te</i>	German	<i>Tee</i>	Greek	τέλον telon
Hebrew	תֵּה, <i>te</i>	Hungarian	<i>tea</i>	Icelandic	<i>te</i>	Indonesian	<i>teh</i>	Irish	<i>tee</i>
Italian	<i>ti</i>	Javanese	ତୀମ୍ବାଳ୍ପୁ	Kannada	ತೀମ୍ବାଳ୍ପୁ	Khmer	ពោត <i>tee</i>	Latin	scientific
Latvian	<i>tēja</i>	Leonese	<i>te</i>	Lithuanian	<i>arbata₍₁₎</i>	Low Saxon	<i>Tee</i> [tɛː] or <i>Tei</i> [taɪ]		<i>thea</i>
Malay	<i>teh</i>	Malayalam	അമ്പാട്ടില	Maltese	<i>te</i>	Norwegian	<i>te</i>	Occitan	<i>te</i>
Polish	herbatka ₍₁₎	Scots	<i>tee</i> [t̬i] – <i>be</i>	Scottish Gaelic	ᚢ. <i>teathu</i>	Sinhalese	තෙය	Spanish	<i>te</i>
Sundanese	entéh	Swedish	<i>te</i>	Tamil	ஏந்தி மென்ற	Telugu	ఎந்தி மென்ற	Western Ukrainian	jerbatka ₍₁₎
Welsh	<i>te</i>			(2)					

Derivatives of *cha* [edit]

Language	Name	Language	Name	Language	Name	Language	Name	Language	Name
Chinese	茶 <i>Chá</i>	Assamese	চাৰ <i>sah</i>	Bengali	চি চা (sai in Eastern regions)	Kapampangan	<i>cha</i>	Cebuano	<i>tsá</i>
English	<i>cha</i> or <i>char</i>	Gujarati	ચા <i>chā</i>	Japanese	茶, ちや <i>cha</i> ₍₁₎	Kannada	ಛಹಾ <i>chahā</i>	Khasi	<i>sha</i>
Punjabi	ਚਾਹੜ ਚਾਹੜ <i>chā</i>	Korean	차 <i>cha</i> ₍₁₎	Kurdish	قا <i>qa</i>	Lao	ຂ້າ /ຂ້າຕົວ <i>Marathi</i>		චਾਹਾ <i>chahā</i>
Odia	ଚାହେଣ୍ଟା <i>cha'a</i>	Persian	چا <i>chā</i>	Portuguese	chá	Sindhi	چاهئن <i>chahen</i>	Somali	شاah <i>shaah</i>
Tagalog	tsaā	Thai	ชา ທ່ອະຊາຍ/	Tibetan	ཇ བྷ/ <i>ja</i>	Vietnamese	trà and <i>chè</i> ₍₂₎		



30

Words -> Word order

Languages typically choose word orders

Word Order	Example	Languages
Subject-Verb-Object (SVO)	I went to school	English, Chinese, Italian
Subject-Object-Verb (SOV)	Mi shaalela gelo (I to-school went)	Japanese, Marathi, Tamil
Verb-Subject-Object (VSO)	Chuaigh mé ar scoil (Went I to-school)	Arabic, Irish
....		

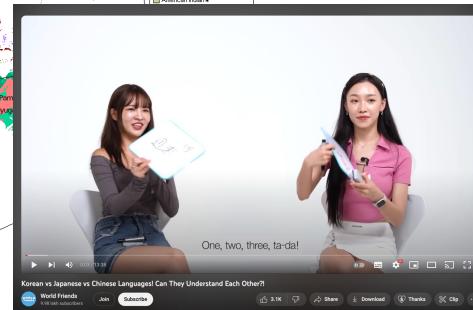
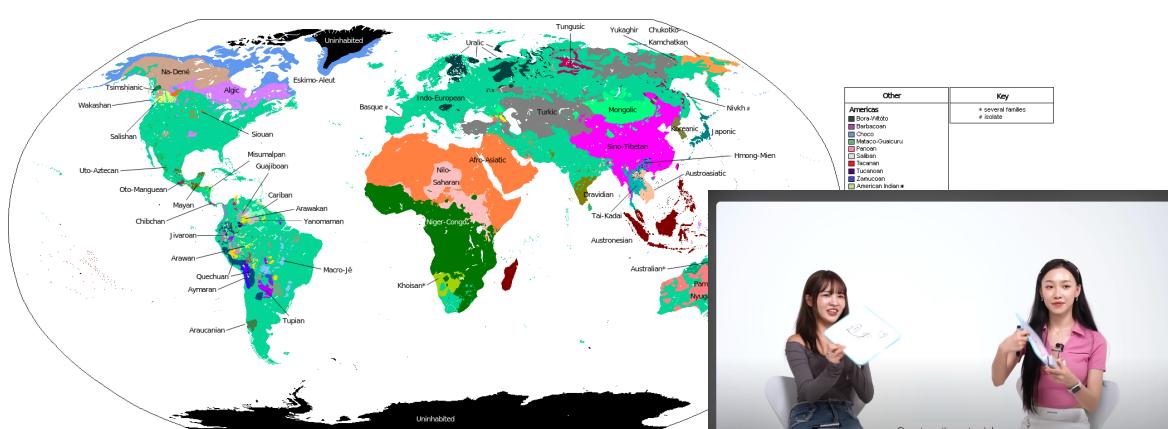
Languages exhibit similarities in terms of word order, word choices, etc.

Tomlin R. S. (1986). *Basic word order : functional principles*. Croom Helm.



31

Language Families



https://en.wikipedia.org/wiki/Language_family

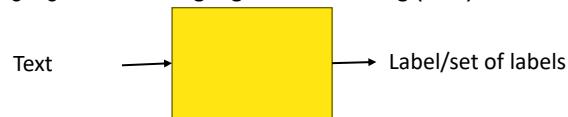
<https://www.youtube.com/watch?v=bYi2pu6gNY0>



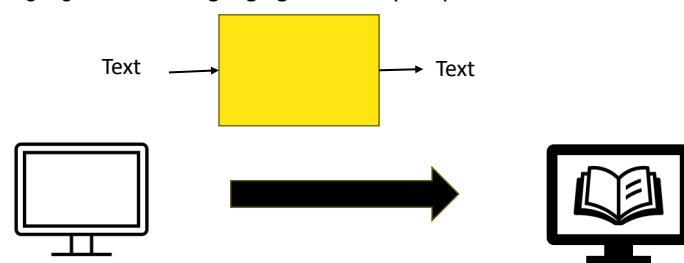
32

So.. What is NLP?

Computers 'understand' language: **Natural language understanding (NLU)**

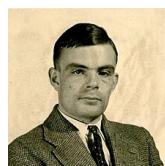


Computers 'generate' language: **Natural language generation (NLG)**

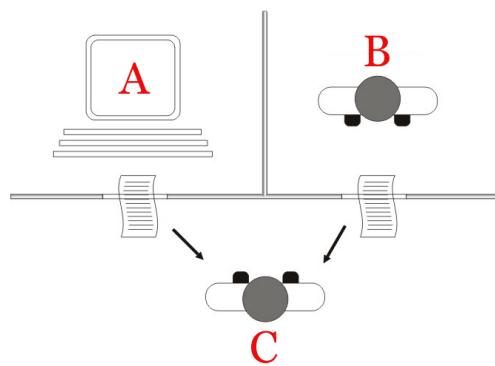


33

The Imitation Game



We now ask the question, 'What will happen when a machine takes the part of A in this game?'



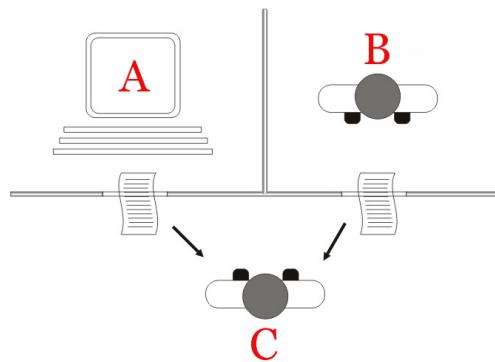
If A fools C into believing that it is human, is it the pinnacle of linguistic machine intelligence? → **NLP**

"**Computing Machinery and Intelligence**", Alan Turing, 1950.



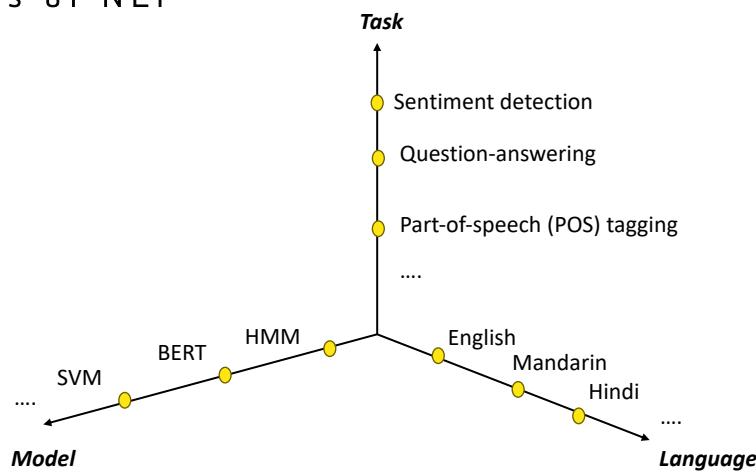
34

What does A need to do to fool C?



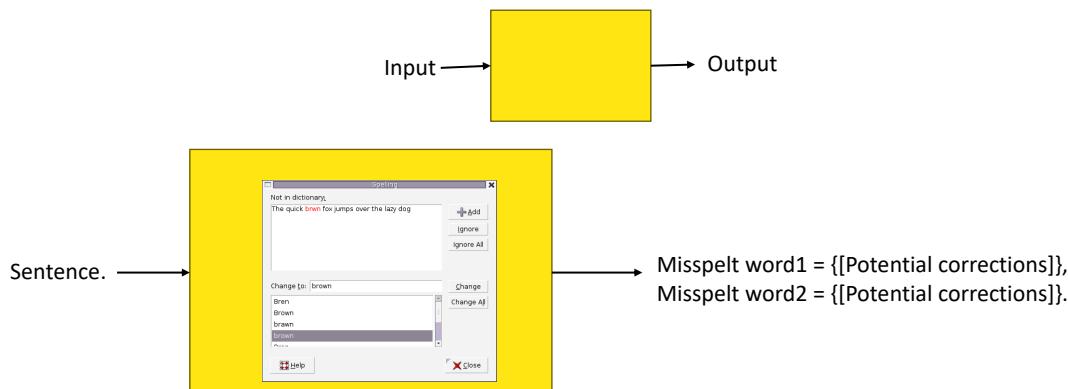
35

Dimensions of NLP



36

Black-box view

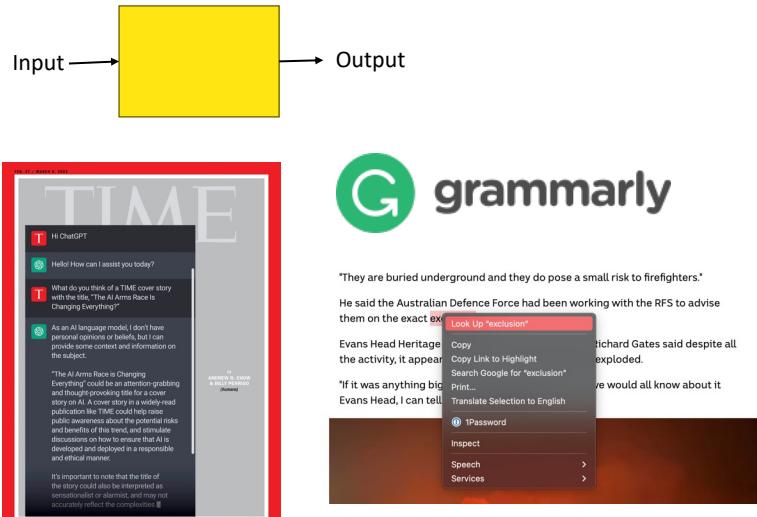


*Any colour that makes you think 'opaque'



37

...what about...?



38

19

Examples of linguistic tasks

Part-of-speech (POS) tagging: Given a sentence, label individual words with POS tags. POS tags are defined by a tag set.

The boy goes to school. → [Yellow Box] → DT NN VB PRP NN .

Chunking: Given a sentence, create chunks by grouping words together. Chunks are typically marked with brackets.

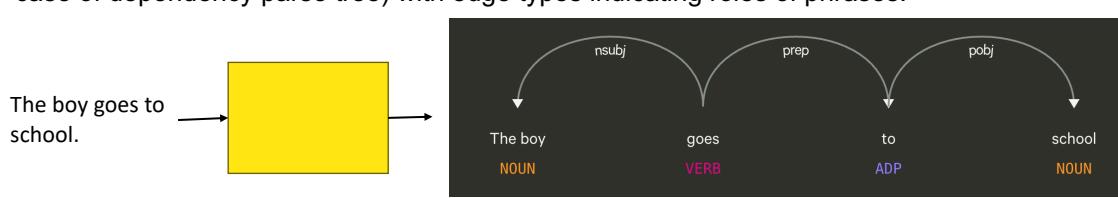
The blue book is on the table. → [Yellow Box] → [The blue book] is on [the table].



39

Examples of linguistic tasks

Parsing: Given a sentence, create a tree rooted at the main verb of the sentence (as in the case of dependency parse tree) with edge types indicating roles of phrases.



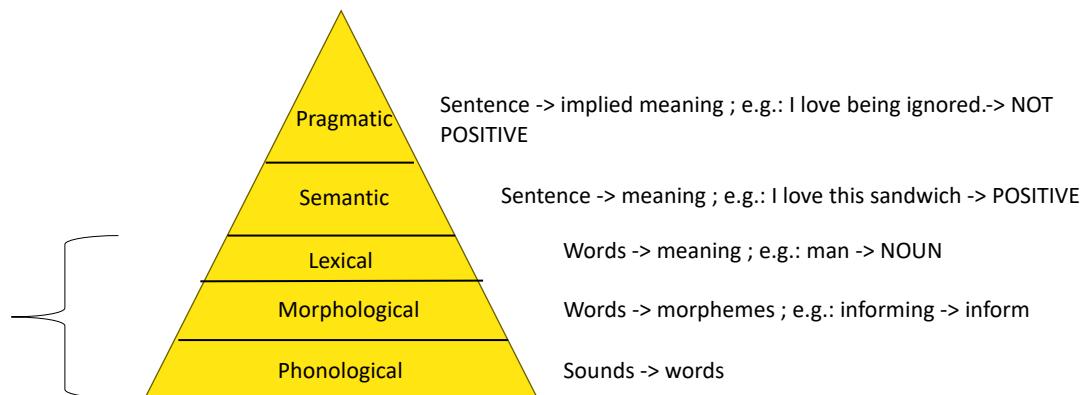
Co-reference resolution: Given a document, link references (such as pronouns) to appropriate phrases in the preceding text.

The blue book is on the table.
Can you pass it to me please? → [Yellow Box] → [The blue book]:1 is on [the table]:2. Can you pass [it]:@1 to me please?



40

These tasks constitute the lower level of the ambiguity hierarchy



41

Demo: NLTK

Natural Language Toolkit (NLTK) (<https://www.nltk.org/>)
 Text processing pipelines for multiple languages
 Several corpora
 Often used as a component in a larger system (say, recommender system)



Demo time!



Bird, Steven, Edward Loper and Ewan Klein (2009), *Natural Language Processing with Python*. O'Reilly Media Inc.



42

Downstream tasks

Downstream tasks: A term used to refer to application-oriented tasks

NLP witnessed a shift in the focus of tasks from linguistic tasks to application-oriented tasks

So what are some downstream tasks?



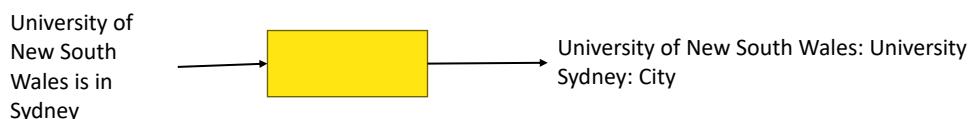
43

Examples of downstream tasks (simplified)

Sentiment Analysis: Given a text, detect the sentiment of the text.



Information Extraction: Given a text, extract entities along with their types. (e.g., named entity recognition)



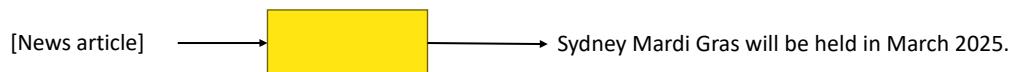
44

Examples of downstream tasks (simplified)

Machine translation: Given a text in a source language, generate its translation in a target language.



Summarisation: Given a document, generate a summary that is shorter than the document but contains crucial information in the document.



45

NLP tasks & research tracks

ACL Anthology: <https://aclanthology.org/>

- Computational Social Science and Cultural Analytics
- Dialogue and Interactive Systems
- Discourse and Pragmatics
- Efficient/Low-Resource Methods for NLP
- Ethics, Bias, and Fairness
- Generation
- Information Extraction
- Information Retrieval and Text Mining
- Interpretability and Analysis of Models for NLP
- Linguistic theories, Cognitive Modeling and Psycholinguistics
- Machine Learning for NLP
- Machine Translation
- Multilingualism and Language Diversity
- Multimodality and Language Grounding to Vision, Robotics and Beyond
- NLP Applications
- Phonology, Morphology and Word Segmentation
- Question Answering
- Resources and Evaluation
- Semantics: Lexical
- Semantics: Sentence-level Semantics, Textual Inference and Other areas
- Sentiment Analysis, Stylistic Analysis, and Argument Mining
- Speech recognition, text-to-speech and spoken language understanding
- Summarization
- Syntax: Tagging, Chunking and Parsing



46

**While the focus of this course will be English, several considerations should get you to think:
“how does this work in languages other than English that I speak?”**

- Data
- Model
- Linguistic phenomena



47



48

Announcements

This week:

No assessment or tutorials this week

Sample quiz for week 1 will be available on Moodle later today – but will not be evaluated. The caption in the quiz clearly says so.

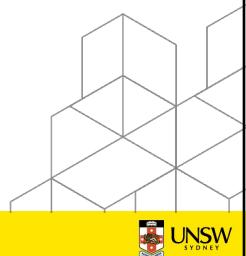
Week 2 & thereafter:

Quizzes from week 2 onwards will be available from Friday evening onwards and will run until Wednesday of next week

Quiz for week 2 will include material from weeks 1 and 2

Consultation Hour: Thursdays 10am-11am | Weeks 2 to 10 | K-17, 217-B.

Start planning your project groups; talk to people ☺ We will ask for group details in week 4



49

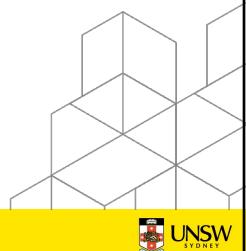
Language conveys ideas

Sydney Opera House



But language is inherently ambiguous, and may convey multiple ideas.

Ambiguity: The quality of being open to more than one interpretation



50

Types of ambiguity

Lexical ambiguity: Words have multiple meanings

"Why should you never date a tennis player?" "Love means nothing to them"

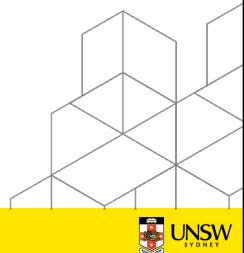
Syntactic ambiguity: Syntactic structures may be ambiguous

"I saw a boy with a telescope."

"I got a job offer from SEEK".

Pragmatic ambiguity: Sentences might imply a different meaning

"Being stranded in traffic is the best way to start the week."

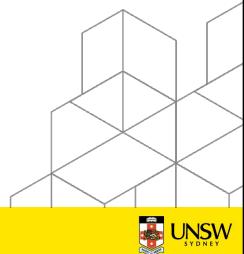
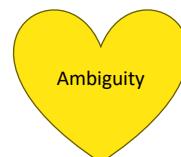


51

Role of Ambiguity in NLP

NLP creates computational techniques to resolve ambiguity in language tasks

Ambiguity resolution in language is at the heart of NLP



52

HELP THE COMPUTER TO FIND THE END OF A SENTENCE

A common task that a computer needs to do with text is to identify the words and the sentences.

Easy for humans

We would like to write a rule to split a document into sentences.

Can you define a rule to be given to a computer?

The Australian Computational and Linguistics Olympiad: <https://ozclo.org.au/>



53

Let's make some rules

The Bank of New York ADR Index, which tracks depository receipts traded on major U.S. stock exchanges, gained 1.3% to 183.32 points in recent session. The index lost 4.63 from the beginning of July. American Depository Receipts are dollar-denominated securities that are traded in the U.S. but represent ownership of shares in a non-U.S. company.

Q. 1. What is the underlying ambiguity?

Q.2. Can you define IF THEN ELSE rules to segment the text into sentences?

Example: "If a character is a '.', start a new sentence"



54

How could have ELIZA done this in Prolog?

```
punctuation(',').
punctuation('.').
punctuation(';').
punctuation('?').
punctuation('!').

sentence_separator(X) :- punctuation(X).

split_to_sentences(Words, Sentences) :-
    split_to_sentences_aux(Words, Sentences, [], []).

split_to_sentences_aux([], Sentences, Sentence_Acc, Sentences_Acc) :-
    ( Sentence_Acc = []
    -> reverse(Sentences_Acc, Sentences)
    ; reverse(Sentence_Acc, Sentence),
      reverse([Sentence | Sentences_Acc], Sentences_Acc)
    ).

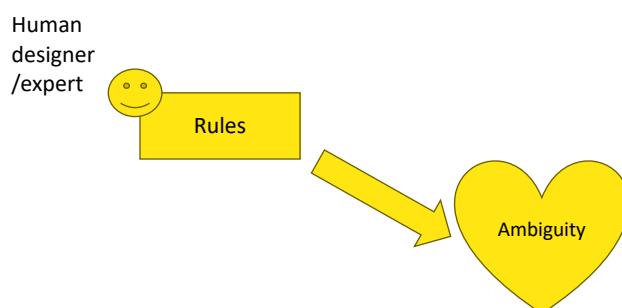
split_to_sentences_aux([X | Xs], Sentences, Sentence_Acc, Sentences_Acc) :-
    ( sentence_separator(X)
    -> reverse(Sentence_Acc, Sentence),
      split_to_sentences_aux(Xs, Sentences, [Sentence | Sentence_Acc], Sentences_Acc)
    ; split_to_sentences_aux(Xs, Sentences, [X | Sentence_Acc], Sentences_Acc)
    ).
```

<https://github.com/bartosz-witkowski/books/tree/master/natural-language-processing-for-prolog-programmers>



55

Rules were an early way to resolve ambiguity in NLP



56

How well do rules work?

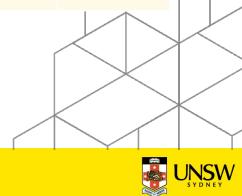
The Bank of New York ADR Index, which tracks depository receipts traded on major U.S. stock exchanges, gained 1.3% to 183.32 points in recent session. The index lost 4.63 from the beginning of July. American Depository Receipts are dollar-denominated securities that are traded in the U.S. but represent ownership of shares in a non-U.S. company.

Precision: % of times the rule works correctly

Recall: % of times the rule covers the positive cases

Rule	#retrieved sentences	#correct sentences	#expected sentences	Precision	Recall
If a character is a ‘’, start a new sentence	12	0	3	0/12 = 0	0/3 = 1
If a character is a ‘’, is followed by a blank space and a capital letter, start a new sentence	3	3	3	3/3 = 1	3/3 = 1

When will the second rule not work? it will not work for this text, for example.



Rules are inadequate to capture linguistic phenomena

Early NLP systems relied on rules

ELIZA: Written in LISP

Rules are time-consuming and expert-intensive

Rules may be subjective and fraught with biases of the ‘expert’

Rules can ‘overfit’

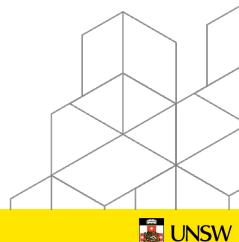
High precision, low recall – typical for rule-based systems

Language is an evolving entity

New words: *lit, kickass*

New concepts: I googled

New events: COVID



Ambiguity can be resolved using observations

Abma is spoken by more than 8,000 people making it one of the largest indigenous languages of Vanuatu, a Pacific island nation that enjoys great linguistic diversity.

Can you translate the following sentences to Abma:

1. The teacher carries the water down.
2. The child keeps eating.
3. The child crawls here.
4. The teacher walks uphill.
5. The palm-tree keeps growing downwards.

NOTE: There is no separate word for 'the' or 'he' in these Abma sentences.

ABMA	ENGLISH
<i>Mwamni sileng.</i>	He drinks water.
<i>Nutsu mwatbo mwamni sileng.</i>	The child keeps drinking water.
<i>Nutsu mwegau.</i>	The child grows.
<i>Nutsu mwatbo mwegalgal.</i>	The child keeps crawling.
<i>Mworob mwabma.</i>	He runs here.
<i>Mwerava Mabontare mwisib.</i>	He pulls Mabontare down.
<i>Mabontare mwisib.</i>	Mabontare goes down.
<i>Mweselkani tela mwesak.</i>	He carries the axe up.
<i>Mwelebte sileng mwabma.</i>	He brings water.
<i>Mabontare mworob mwesak.</i>	Mabontare runs up.
<i>Sileng mworob.</i>	The water runs.

Now here are some more words in Abma:

ABMA	ENGLISH
<i>sesesrakan</i>	teacher
<i>mwiegani</i>	eat
<i>buet</i>	taro
<i>muhural</i>	walk
<i>butsu-kul</i>	palm-tree

<https://ozclo.org.au/wp-content/uploads/2014/11/2010-first-round-problems.pdf>



59

Solution

ENGLISH	ABMA
1. The teacher carries the water down.	<i>Sesesrakan mweselkani sileng mwisib.</i>
2. The child keeps eating.	<i>Nutsu mwatbo mwiegani.</i>
4. The child crawls here.	<i>Nutsu mwegalgal mwabma.</i>
5. The teacher walks uphill.	<i>Sesesrakan muhural mwesak.</i>
6. The palm-tree keeps growing downwards.	<i>Butsukul mwatbo mwiegau mwisib.</i>



60

30

Let's introspect...

How did you use the data?

What assumptions did you make?

Challenge: What sentences would not be possible to translate with certainty?

NOTE: There is no separate word for 'the' or 'he' in these Abma sentences.

ABMA	ENGLISH
<i>Mwamni sileng.</i>	He drinks water.
<i>Nutsu mwatbo mwamni sileng.</i>	The child keeps drinking water.
<i>Nutsu mwegau.</i>	The child grows.
<i>Nutsu mwatbo mwegalgal.</i>	The child keeps crawling.
<i>Mvorob mwabma.</i>	He runs here.
<i>Mwerava Mabontare mwisib.</i>	He pulls Mabontare down.
<i>Mabontare mwisib.</i>	Mabontare goes down.
<i>Mweselkani tela mwesak.</i>	He carries the axe up.
<i>Mwelebte sileng mwahma.</i>	He brings water.
<i>Mabontare mavorob mwesak.</i>	Mabontare runs up.
<i>Sileng mworob.</i>	The water runs.

Now here are some more words in Abma:

ABMA	ENGLISH
<i>sesesrakan</i>	teacher
<i>mwegani</i>	eat
<i>buet</i>	taro
<i>muhural</i>	walk
<i>butsu-kul</i>	palm-tree



61

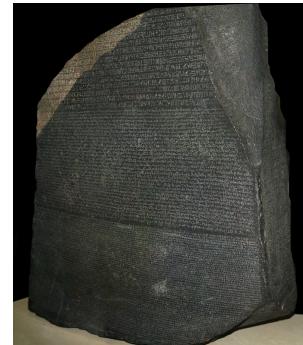
Observations come from data

Rosetta stone

196 BC

Three versions of a decree

One of the earliest known 'corpus'

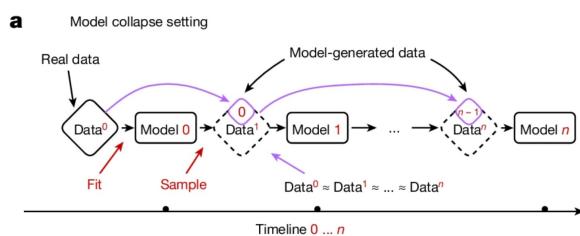


Textual Corpus/Dataset

Digitisation made datasets available

Internet exploded available datasets

Website have been blocking
crawlers that collect data



<https://www.theguardian.com/technology/2023/aug/25/new-york-times-cnn-and-abc-block-openai-gptbot-web-crawler-from-scraping-content>
<https://www.nature.com/articles/s41586-024-07566-y>



62

Terminology

Supervised:

Data annotated with expected output labels
 "I love the movie" -> Positive ; "The movie sucks" -> Negative

Distant-supervised:

Data annotated with expected output labels, using heuristics. Labeling is 'distant'
 "I love the movie #not" -> Sarcastic ; 'I love the movie' -> Not sarcastic

Unsupervised:

Data with no annotations w.r.t. the output labels
 "I love the movie"

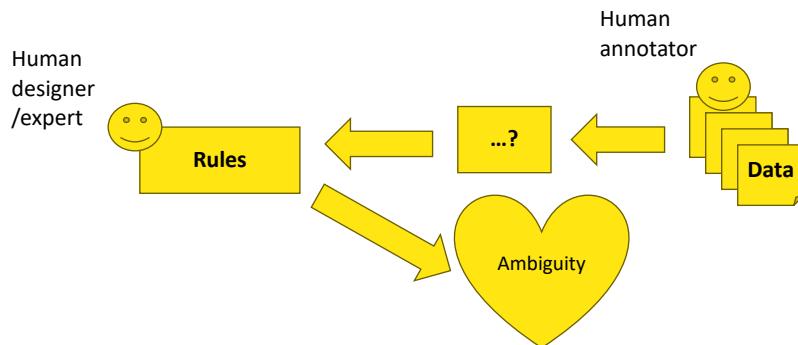
Self-supervised:

The data itself becomes supervision.
 Introduced in deep learning-based NLP
 "I love the movie" is mapped to: "I ____ the movie" -> "love"

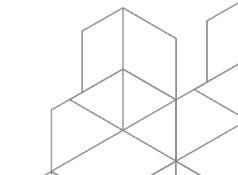


63

Data can help learn 'rules'



... errm... What does learning mean?



64

Enter: Probability




Pen-and-paper time!




65

Re~~c~~ap: Terminology

- Naïve Bayes
- Conditional independence assumption
- N-grams
- 'Skip'-grams?
- Argmax
- Suggested reading: <https://web.stanford.edu/~jurafsky/slp3/4.pdf>
- N-gram language models (Module 2)




66

Are words really random variables?

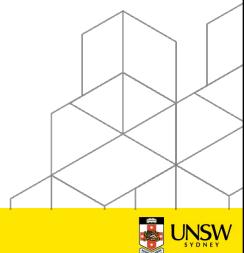
"Her face **fell** when she heard that she had been fired."

"The fruit **fell** from the tree."

"Avoid vulgar language please."

"Now that's real crude behavior!"

One-hot vector -> Word vector representations (Module 2)



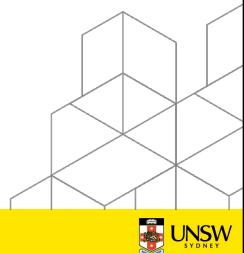
67

softmax

Neural networks generate vectors with logits

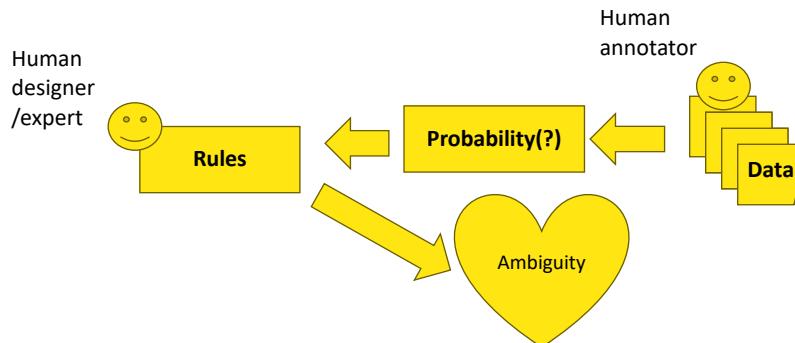
To convert them to probabilities, softmax function is used

It has desirable properties not seen in standard normalisation



68

Data, Computation (Probability) and Ambiguity

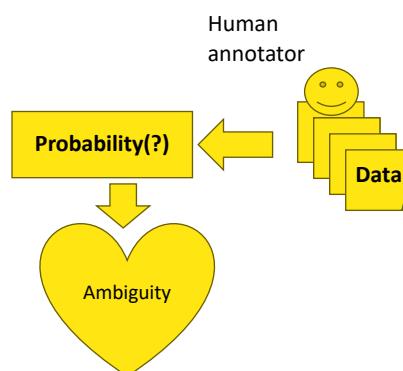


... errm.. What does learning mean?



69

Data, Computation (Probability) and Ambiguity



70

Can 'data' be anything other than textual data? Lexicons

WordNet (<https://wordnet.princeton.edu/>)

LIWC (Lexical Inquiry and Word Count)

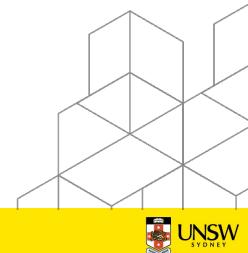
Domain-specific lexicons

Medical ontologies

Proprietary assets: Businesses often consider in-house lexicons to be valuable assets

May be represented either as graphs or lists

What is the 'easiest' way to use such lexicons?



71

WordNet

A lexical database

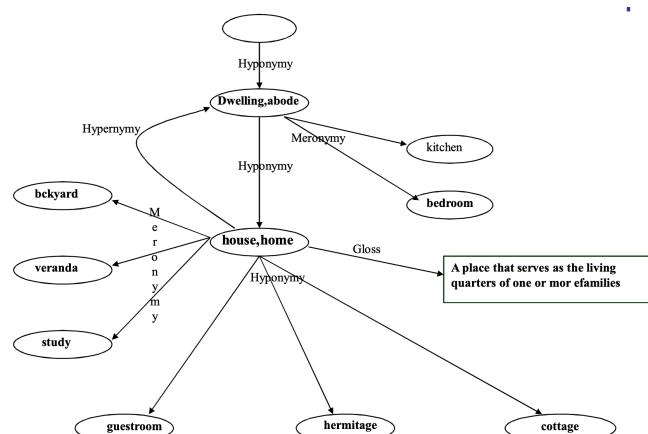
Work on the English language WordNet began in 1985

Inspired WordNets in other languages

<http://wordnetweb.princeton.edu/perl/webwn>

Synsets (synonym sets) connected to each other by relations

Semantic relations: Synonymy, Antonymy, Hypo/hypernymy, mero/holonymy.



Fellbaum, Christiane, ed. *WordNet: An electronic lexical database*. MIT press, 1998.
Example from: <https://www.cse.iitb.ac.in/~pb/cs626-2013/cs626-lect24to26-wn-2013-9-17.pdf>



72

spaCy

spaCy is an open-source software library for NLP

Multi-lingual support

"Industrial-strength NLP"

<https://spacy.io/usage/spacy-101>

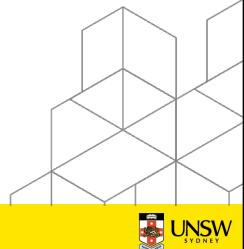
Supports deep learning pipelines

AllenNLP: Built on the top of spaCy and PyTorch (<https://spacy.io/universe/project/allennlp>)

Integration with LLMs available as well

<https://spacy.io/usage/large-language-models>

```
pip install spacy
python -m spacy download en_core_web_sm
```



73

Spacy Matcher

Regular expressions: Another artifact of rule-based NLP

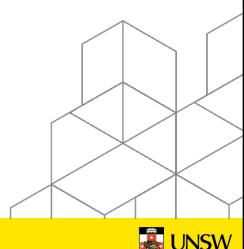
SpaCy matcher allows simple pattern-matching using regular expressions

Several types of matcher primitives: <https://spacy.io/api/matcher>

PhraseMatcher may be useful for ontology-based matching



Demo time!



74

```
doc = self.nlp(text)
print("Sentences are:")
for sent in doc.sents:
    print(sent.text)
```

```
(base) z35399580L-Q601F00993 week1 % python spacy2_sentencesplitter.py "The Bank of New York ADR Index, which tracks depositary receipts traded on major U.S. stock exchanges, gained 1.3% to 183.32 points in recent session. The index lost 4.63 from the beginning of July. American Depository Receipts are dollar-denominated securities that are traded in the U.S. but represent ownership of shares in a non-U.S. company."
```

```
Sentences are:
The Bank of New York ADR Index, which tracks depositary receipts traded on major U.S. stock exchanges, gained 1.3% to 183.32 points in recent session.
The index lost 4.63 from the beginning of July.
American Depository Receipts are dollar-denominated securities that are traded in the U.S. but represent ownership of shares in a non-U.S. company.
```

<https://spacy.io/usage/rule-based-matching>



75

spaCy: Recap

Black-box NLP library for several NLP tasks
SpaCy matcher!



76



77

Computational techniques in NLP have witnessed paradigm shifts

Generation 1: Rule-based NLP

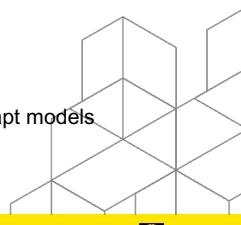
- Human-written rules
- Highly explainable
- Any adaptation needs human intervention

Generation 2: Statistical NLP

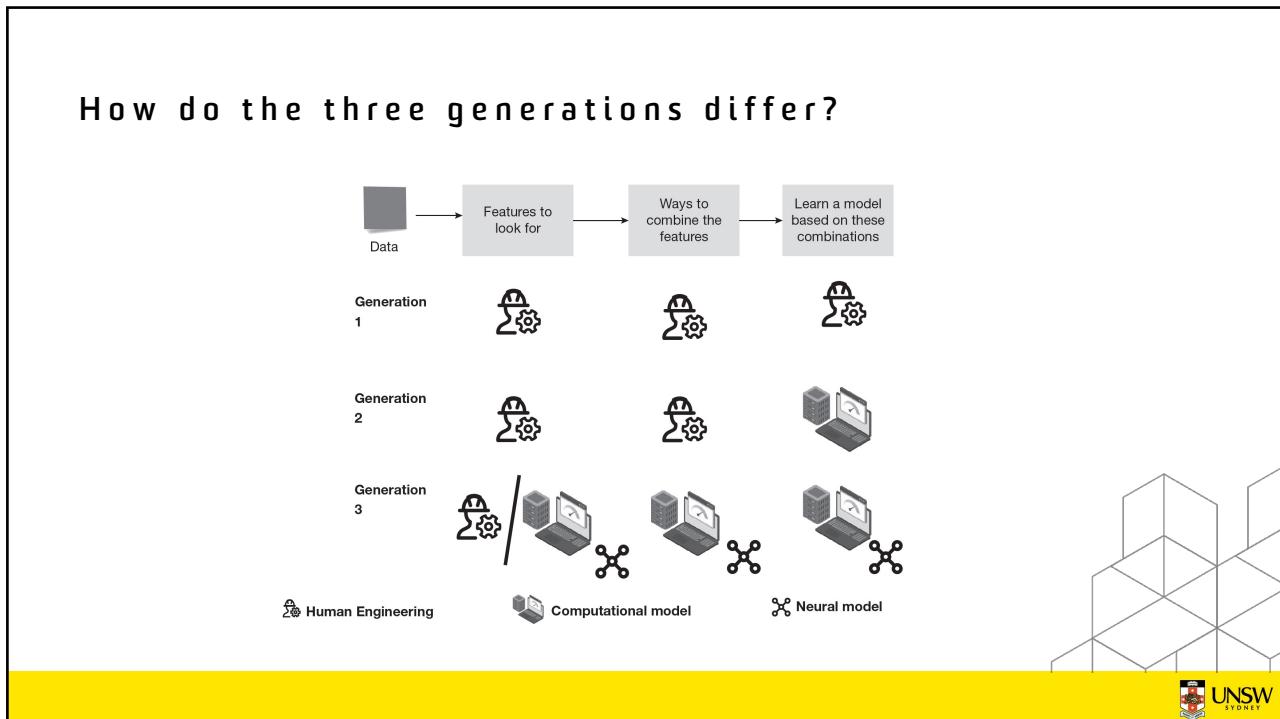
- Human-written features
- Probability and argmax at the heart of computation

Generation 3: Neural NLP

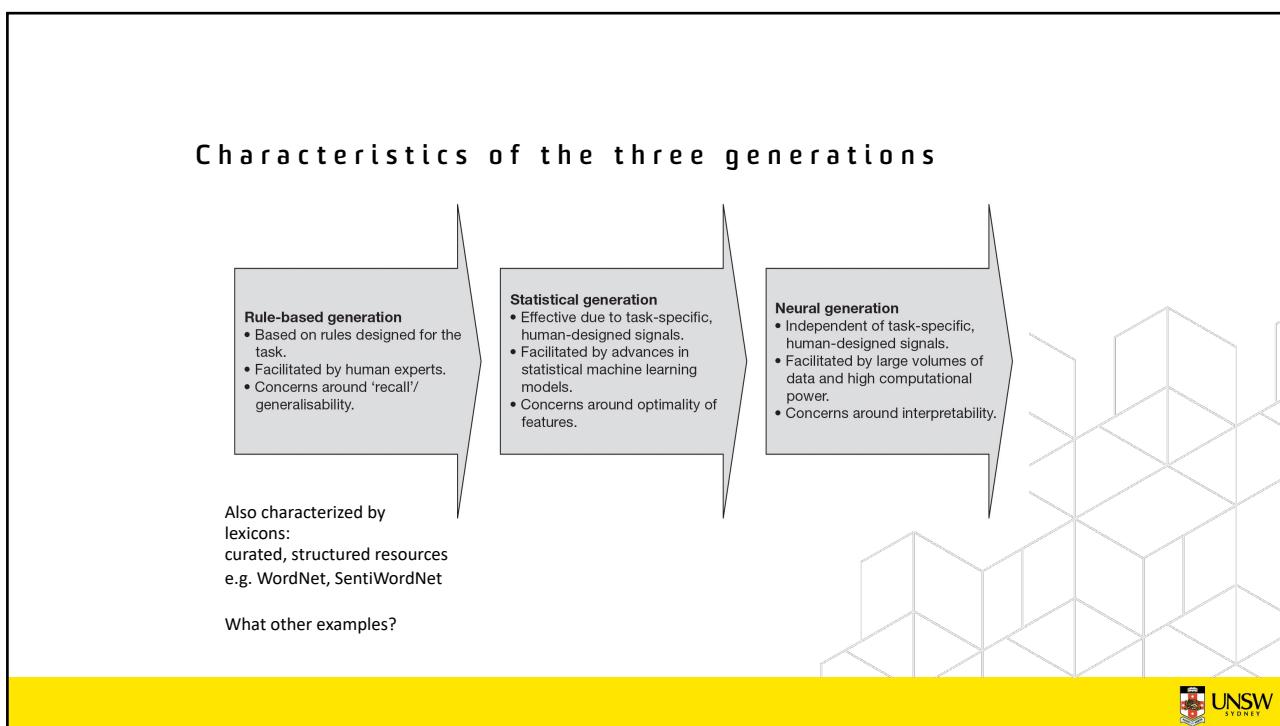
- Overt task-specific features not essential
- Foundation models learn from self-supervised data; Task-specific data may be used to adapt models



78



79



80

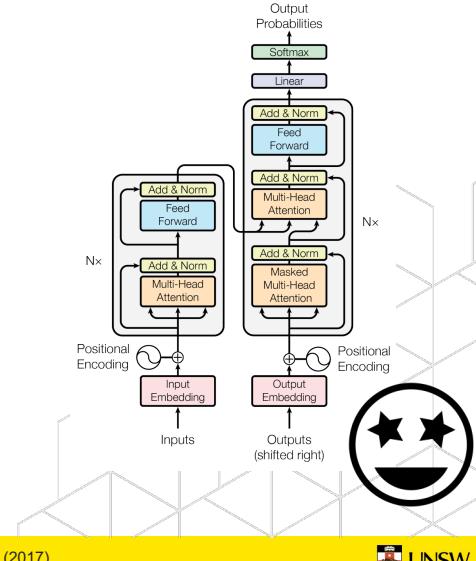
Neural NLP has been influenced by Transformer

Transformer is an architecture originally introduced for NLP; adapted to other areas of AI

Models trained using Transformer architecture

Input Sequence → **Transformer** → Output Sequence

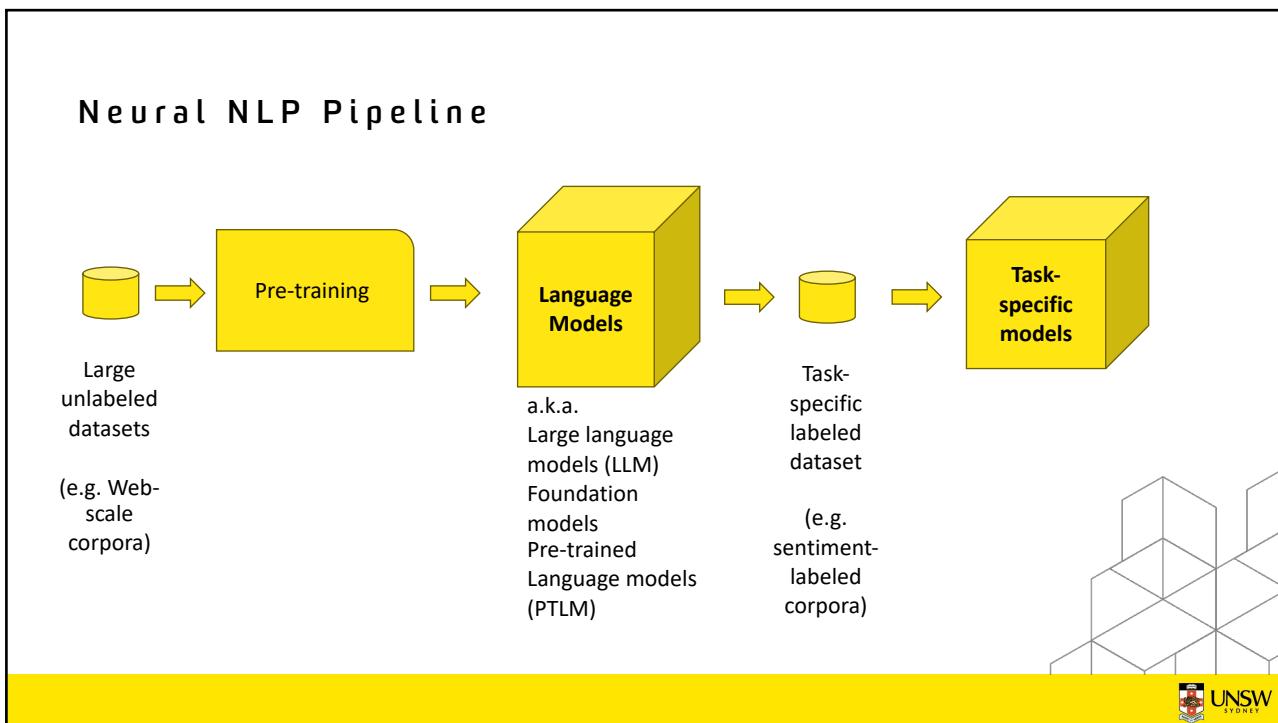
Encoder-only models: **BERT**
Decoder-only models: **GPT** (Week 4)



Vaswani, Ashish, et al. "Attention is all you need." *Advances in neural information processing systems* 30 (2017).

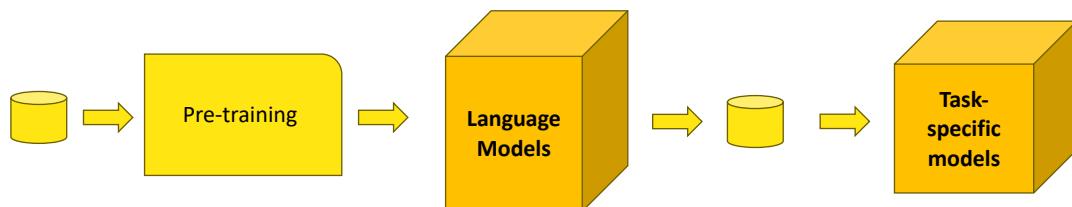
UNSW SYDNEY

81



82

Neural NLP Pipeline



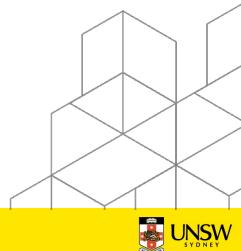
Both PTLMs and TLMs can be used for inference.

How are language models ‘pre-trained’?

How are they adapted to specific tasks?

.. And related considerations.

All in Week 4



83

Another black-box library: HuggingFace

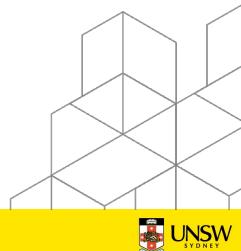
HuggingFace (<https://huggingface.co/>)

HuggingFace makes pre-trained and/or fine-tuned models available

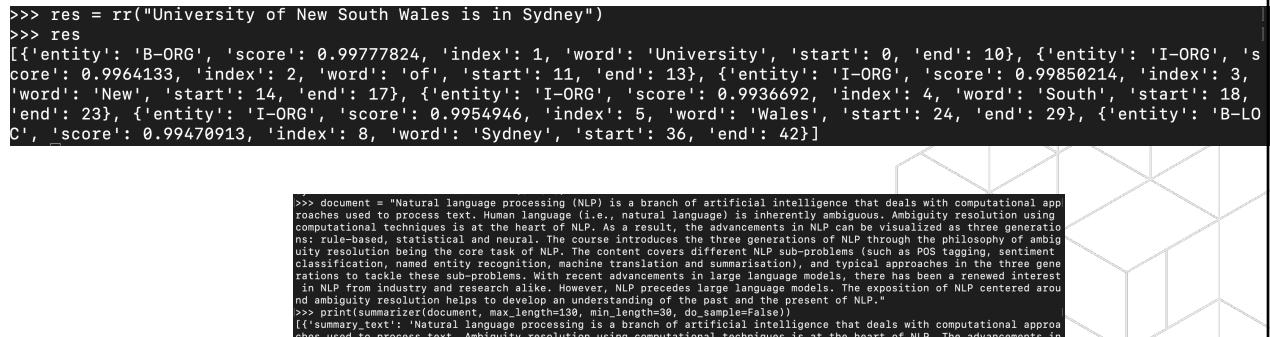
<https://huggingface.co/meta-llama/Llama-2-7b-chat-hf>



Demo time!



84



```
>>> res = rr("University of New South Wales is in Sydney")
>>> res
[{'entity': 'B-ORG', 'score': 0.99777824, 'index': 1, 'word': 'University', 'start': 0, 'end': 10}, {'entity': 'I-ORG', 'score': 0.9964133, 'index': 2, 'word': 'of', 'start': 11, 'end': 13}, {'entity': 'I-ORG', 'score': 0.99850214, 'index': 3, 'word': 'New', 'start': 14, 'end': 17}, {'entity': 'I-ORG', 'score': 0.9936692, 'index': 4, 'word': 'South', 'start': 18, 'end': 23}, {'entity': 'I-ORG', 'score': 0.9954946, 'index': 5, 'word': 'Wales', 'start': 24, 'end': 29}, {'entity': 'B-LOC', 'score': 0.99470913, 'index': 8, 'word': 'Sydney', 'start': 36, 'end': 42}]
```

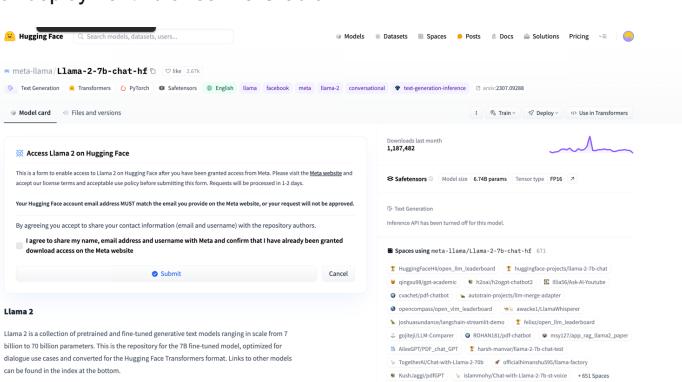
```
>>> document = "Natural language processing (NLP) is a branch of artificial intelligence that deals with computational approaches used to process text. Human language (i.e., natural language) is inherently ambiguous. Ambiguity resolution using computational techniques is at the heart of NLP. As a result, the advancements in NLP can be visualized as three generations: rule-based, statistical and neural. The course introduces the three generations of NLP through the philosophy of ambiguity resolution being the core task of NLP. The content covers different NLP sub-problems (such as POS tagging, sentiment classification, named entity recognition, machine translation, summarization, etc.) and typical approaches in the three generations to tackle these sub-problems. With recent developments in large language models, there has been a renewed interest in NLP from industry and research alike. However, NLP precedes large language models. The exposition of NLP centered around ambiguity resolution helps to develop an understanding of the past and the present of NLP."
>>> print(summarizer(document, max_length=130, min_length=30, do_sample=False))
[{'summary_text': 'Natural language processing is a branch of artificial intelligence that deals with computational approaches used to process text. Ambiguity resolution using computational techniques is at the heart of NLP. The advancements in NLP can be visualized as three generations: rule-based, statistical and neural.'}]
```

UNSW
SYDNEY

85

Hugging Face Hub

Community-contributed repository of models and datasets
Easy to integrate with quick deployment libraries like Gradio



Llama 2

Llama 2 is a collection of pretrained and fine-tuned generative text models ranging in scale from 7 billion to 70 billion parameters. This is the repository for the 7B fine-tuned model, optimized for dialogue use cases and converted for the Hugging Face Transformers format. Links to other models can be found in the index at the bottom.

Model Details

Model ID: llama-2-7b-chat-hf

Model size: 6.740 params

Tensor type: FP16

Downloads last month: 3,387,402

Safetensors: Model size: 6.740 params Tensor type: FP16

Text Generation: Inference API has been turned off for this model.

Spaces using meta-llama/Llama-2-7b-chat-hf: 671

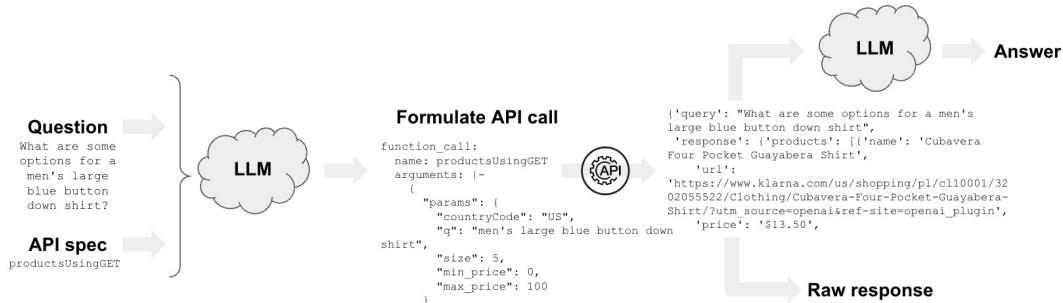
- HuggingFaceProjects_llm_leaderboard
- spqrullt4gt-academy
- coacherJdHtbd
- opennmt-project/mine-megatron
- pythia-100m-100b-1000b
- AlfredLPTF5-chat-GPT
- TogetherLChat-with-Llama-3-11B
- KoHuggJyaffF
- IsarandyJChat-with-Llama-2-7b-or-10
- 462 Spaces

UNSW
SYDNEY

86

APIs for LLMs

Commercial LLMs offer access in the form of API calls.
Often in the form of an API call.



https://python.langchain.com/v0.1/docs/use_cases/apis/



87

Example: OpenAI API



Example: OpenAI API call

```
from openai import OpenAI
client = OpenAI()

response = client.chat.completions.create(
    model="gpt-4o",
    messages=[
        {
            "role": "system",
            "content": "You will be provided with statements, and your task is to convert them to standard English."
        },
        {
            "role": "user",
            "content": "She no went to the market."
        }
    ],
    temperature=1,
    max_tokens=256,
    top_p=1
)
```

https://python.langchain.com/v0.1/docs/use_cases/apis/



88

... but ...

The screenshot shows a Reddit post on the r/Startup_Ideas subreddit. The post is titled "Thoughts on llm based startups". The content discusses the observation that many AI startups are building wrappers on GPT or other LLMs for functionalities that are already capable of doing on their own like chatbot, summarization, etc. It notes that big players are bringing all AI usecases known in public domain in their products. The post also mentions that most AI startups are not doing anything fundamentally different than what GPT can already do, and many are just inventing usecases to solve it using AI somehow. A yellow bar at the bottom right contains the UNSW Sydney logo.

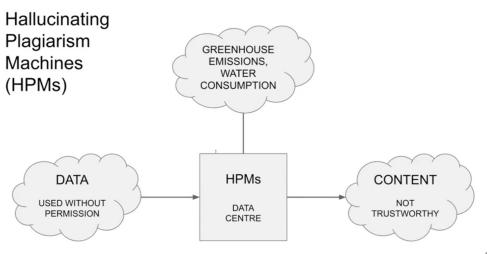
89

The diagram features a yellow square in the top left corner containing the UNSW Sydney logo and the text "Australia's Global University". The main area consists of a complex network of light gray lines forming a three-dimensional geometric structure. In the center, the text "Part 4" is written above "Key Considerations" in a bold, sans-serif font. A horizontal line is positioned below "Key Considerations".

90

Is NLP a 'solved' problem?

Availability of web-based demo may give the impression that NLP is a 'solved' problem



Several considerations assume importance.

I hope this course will encourage you to investigate some of the above aspects.

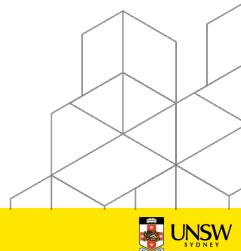


Image from Steven Bird's ALTA 2024 keynote: https://www.youtube.com/watch?v=YcU_tZlesKI

91

Hallucination

In the case of generative models

Factually incorrect

Semantically inconsistent

.. Other definitions of hallucination are possible.

Prompt: Engineering effort to build Eiffel tower ALARMING

AI-generated text: ...Designed by Gustave Eiffel, it was inaugurated in 1889 to celebrate the 100th anniversary of the European Civil War...

Fact: Eiffel tower was built to celebrate the 100th anniversary of the French Revolution.

Which LLM might work well for you?

Chatbot Arena: <https://lmarena.ai/>

*Rawte, Vipula, et al. "The Troubling Emergence of Hallucination in Large Language Models—An Extensive Definition, Quantification, and Prescriptive Remediations." EMNLP 2023.



92

Ethical Considerations

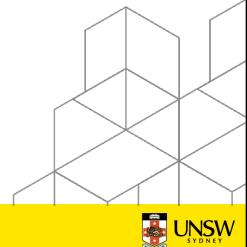
Bias

Sentence completion can produce sexist/racist/religiously offensive output

Privacy

Transparency

Explanation of output; citing one's source



93

Gender Bias in modern NLP

Task	Example of Representation Bias in the Context of Gender	D	S	R	U
Machine Translation	Translating “He is a nurse. She is a doctor.” to Hungarian and back to English results in “She is a nurse. He is a doctor.” (Douglas, 2017)		✓	✓	
Caption Generation	An image captioning model incorrectly predicts the agent to be male because there is a computer nearby (Burns et al., 2018).		✓	✓	
Speech Recognition	Automatic speech detection works better with male voices than female voices (Tatman, 2017).			✓	✓
Sentiment Analysis	Sentiment Analysis Systems rank sentences containing female noun phrases to be indicative of anger more often than sentences containing male noun phrases (Park et al., 2018).		✓		
Language Model	“He is doctor” has a higher conditional likelihood than “She is doctor” (Lu et al., 2018).		✓	✓	✓
Word Embedding	Analogy such as “man : woman :: computer programmer : homemaker” are automatically generated by models trained on biased word embeddings (Bolukbasi et al., 2016).	✓	✓	✓	✓

Table 1: Following the talk by Crawford (2017), we categorize representation bias in NLP tasks into the following four categories: (D)enigration, (S)tereotyping, (R)eognition, (U)nder-representation.

*Tony Sun, Andrew Gaut, Shirlyn Tang, Yuxin Huang, Mai ElSherief, Jieyu Zhao, Diba Mirza, Elizabeth Belding, Kai-Wei Chang, and William Yang Wang. 2019. Mitigating Gender Bias in Natural Language Processing: Literature Review. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 1630–1640, Florence, Italy. Association for Computational Linguistics.



94

Dialectal Bias in modern NLP

NLP Task	Paper	Impact
Language classification	[Blodgett et al. 2016]	Language detection shows lower performance for African-American English.
Sentiment classification	[Okpala et al. 2022]	Text in African-American English may be predicted higher as hate speech.
Natural Language Understanding	[Ziems et al. 2022]	Popular models perform worse on GLUE tasks for African-American English text.
Summarisation	[Keswani and Celis 2021]	Generated multi-document summaries may be biased towards majority dialect.
Machine translation	[Kantharuban et al. 2023]	Significant drop in MT from and to dialects of Portuguese/Bengali/etc. to and from English.
Parsing	[Scannell 2020]	Lower performance of parsers on Manx Gaelic as compared to Irish/Scottish Gaelic.

Table 1. Examples of adverse impact on NLP task performance due to dialectic variations.

Joshi, Aditya, Raj Dabre, Diptesh Kanodia, Zhuang Li, Haolan Zhan, Gholamreza Haffari, and Doris Dippold. "Natural Language Processing for Dialects of a Language: A Survey." *ACM Computing Surveys* 2025.



95

Privacy in modern NLP

Samsung bans Generative AI use
Data leakage into training and test set

<https://www.bloomberg.com/news/articles/2023-05-02/samsung-bans-chatgpt-and-other-generative-ai-use-by-staff-after-leak>

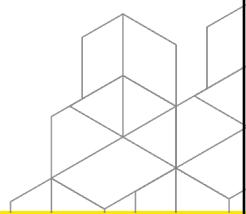


96

Transparency

<https://www.perplexity.ai/>

The screenshot shows the Perplexity AI interface. The main query is "Why is Italian food so popular?". Below the query, there's a section titled "Sources" with links to "Public domain Entry", "Italian Cuisine", "why is it so probably popular?", "Why is Italian Food so Popular?", "publicdomaines... - 1", "italiancuisine... - 2", and "villaino... - 3". There's also a "View 2 more" link. The "Answer" section discusses the popularity of Italian food due to its variety, quality, simplicity, and use of fresh ingredients. It mentions the influence of Italian immigrants and their role in promoting their food. The "Related" section includes links to "what are some popular Italian dishes" and "how has Italian food influenced other countries". On the right side, there are buttons for "Share", "Rewrite", "Search Videos", and "Generate Image". A watermark for "perplexity.ai" is visible at the bottom right.



<https://techcrunch.com/2024/01/04/ai-powered-search-engine-perplexity-ai-now-valued-at-520m-raises-70m/>



97

Summary

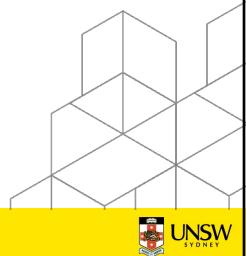
	Key Idea	Demos
NLP Today & Yesterday	NLP has fascinated computer science for a long time.	Fundamental NLP tasks using NLTK
Ambiguity	Interaction between data, probability and ambiguity resolution	Text matching using spaCy
The three generations	Three-generational view of NLP	Open-source NLP models using HuggingFace
Considerations	NLP is far from solved. Hallucination, privacy, biases, etc.	Emerging NLP tools: Perplexity.ai



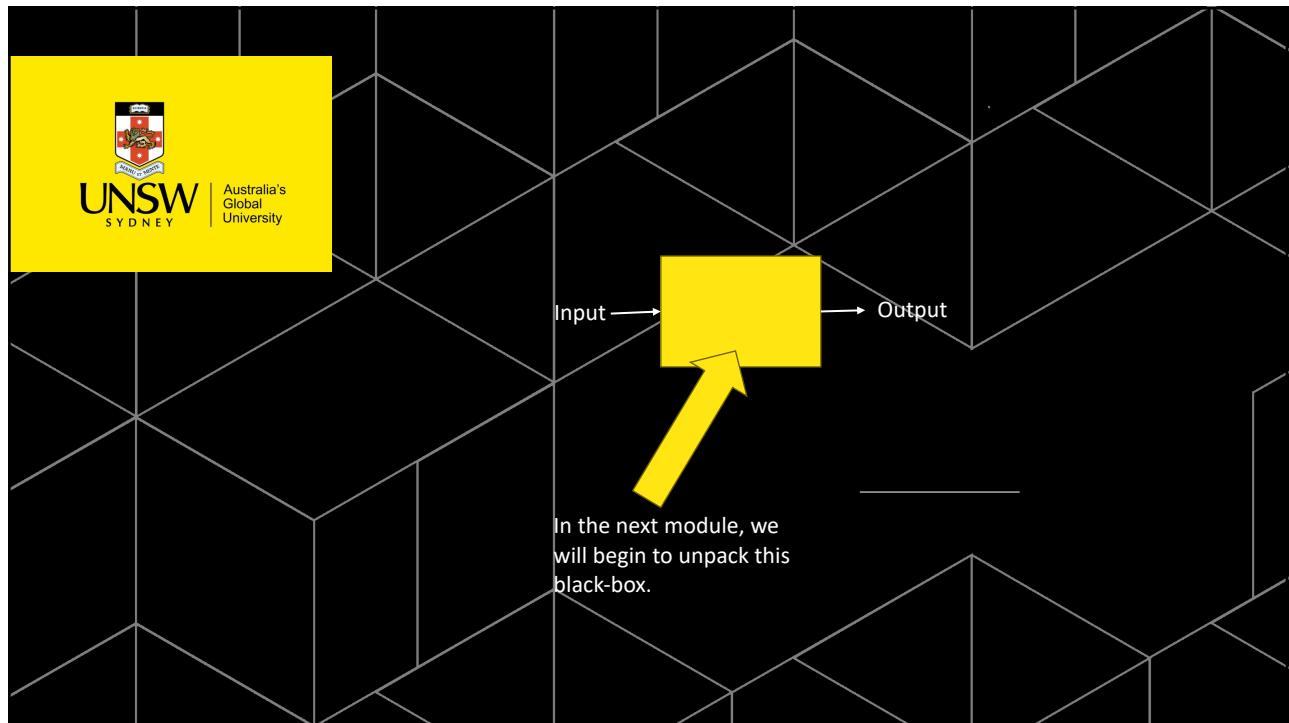
98

Coming up next...

Grammar, Probabilistic language models
Word vectors
Sequential networks



99



100