

# Natural Language Processing

# Week1: Introduction

# Dr Haithem Afli

@AfliHaithem

Haithem.afli@mtu.ie



# Roadmap

- Introduction
  - **Lecturer & module**
  - Topics
  - Syllabus discussion
  - Books and Resources
- Lecture 1

# Lecturer - Haithem Afli

## Computer Science Lecturer at MTU

- NLP, Applied Machine Learning
- MSc in AI, Data Analytics and Computational Biology



## SFI Funded Investigator at ADAPT Centre

- ADAPT lead at MTU
- Member of the ADAPT Executive Committee and challenge lead - Analysing Digital Content
- Manager of SIGMA Research Group



## Research Interest:

- Natural Language Processing
- Social Media and UGC Analysis
- Machine Translation
- Data Analytics and Computational Biology



<https://www.adaptcentre.ie/experts/haithem-afli/>

## Teaching Experience

- 4 years – CIT-MTU (NLP, ML, DA)
- 3 years - DCU (MT, ML)
- 5 years – Le Mans University, France (C/C++Prog., databases)



[www.mtu.ie](http://www.mtu.ie)

Succeeding Together



Horizon 2020  
European Union Funding  
for Research & Innovation

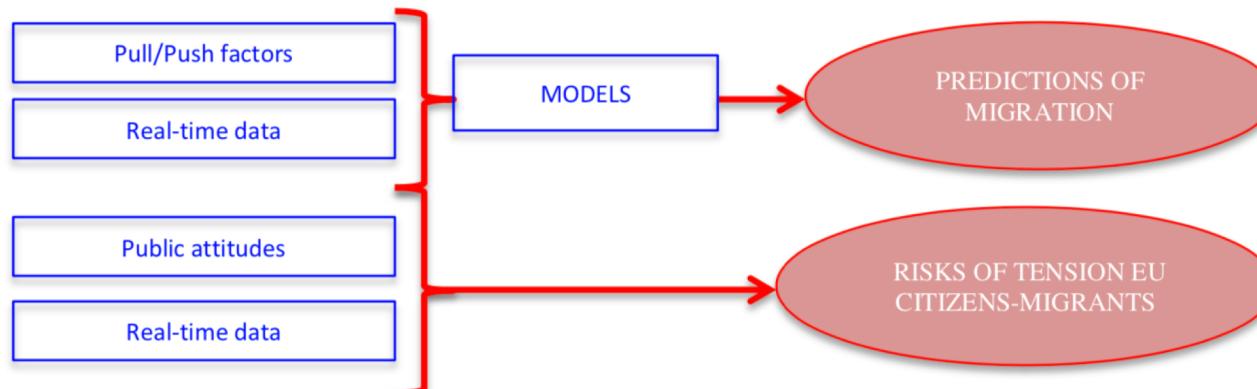


**MTU**

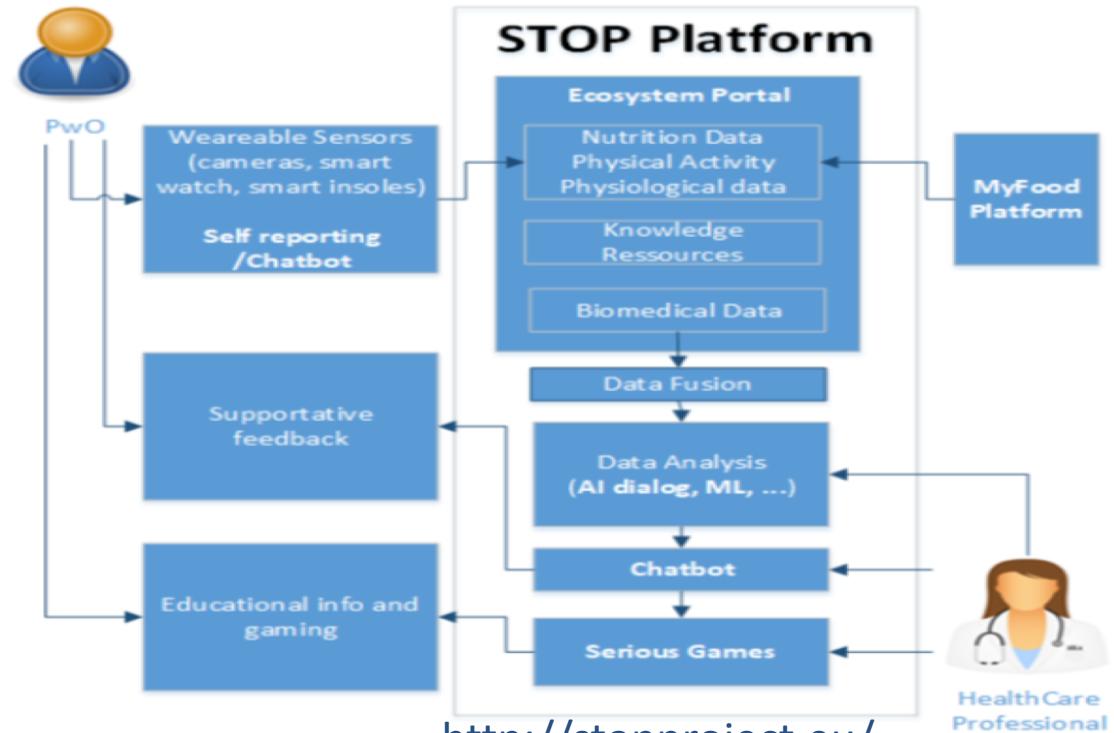
Ollscoil Teicneolaíochta na Mumhan  
Munster Technological University

# itflows

IT TOOLS AND METHODS  
FOR MANAGING  
MIGRATION FLOWS



[www.itflows.eu](http://www.itflows.eu)



<http://stopproject.eu/>

Succeeding Together

[www.mtu.ie](http://www.mtu.ie)

# Labs and support

- Praveen Joshi (for the full-time group)
  - Casual Lecturer at the Department of Computer Science
  - Email: [Praveen.joshi@mycit.ie](mailto:Praveen.joshi@mycit.ie)
- Qualification:
  - Masters in Artificial Intelligence – CIT, 2018-2019
  - PhD candidate in AI
- Projects:
  - Slice Net
- Industrial Exp:
  - Infosys
  - Siemens
  - Accenture AI
  - Clear stream
  - AIP
  - Speire



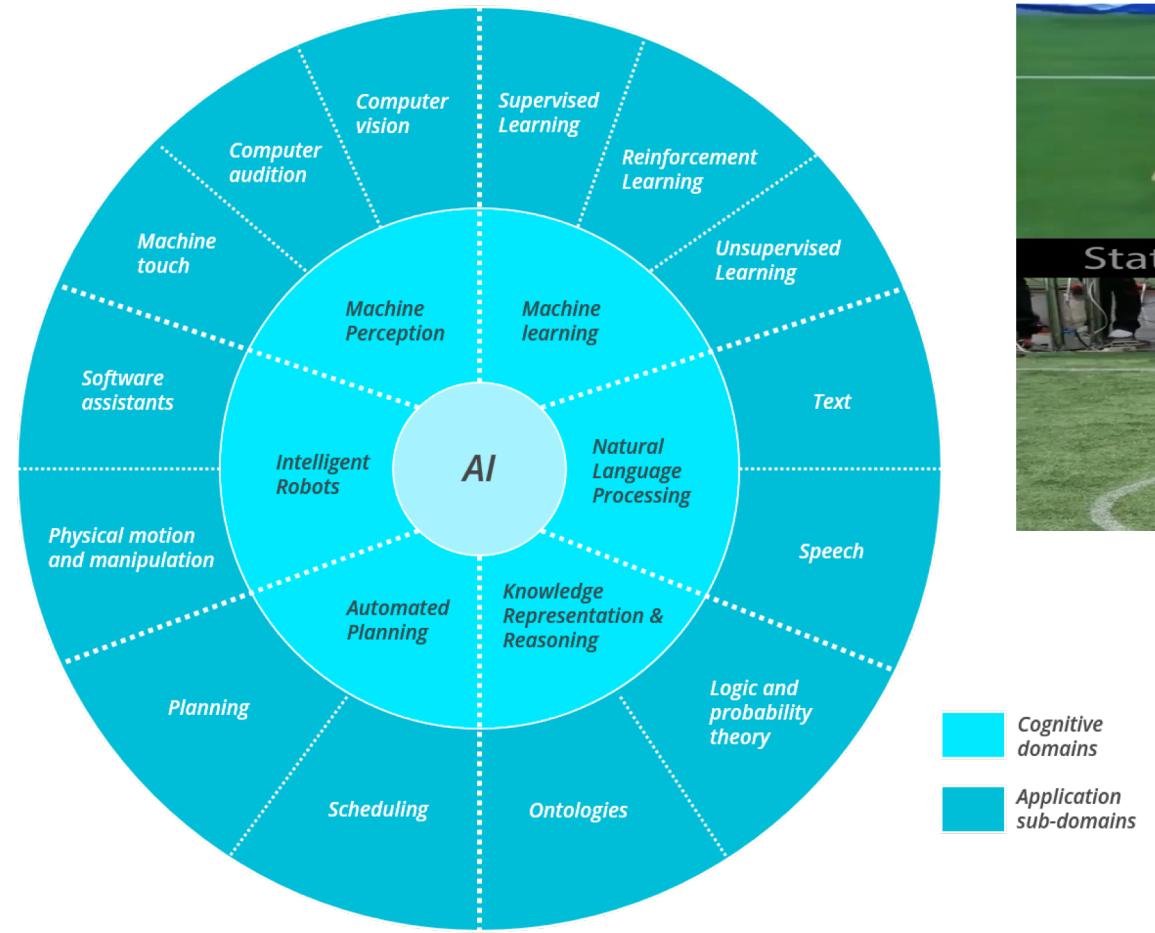
**MTU**

Ollscoil Teicneolaíochta na Mumhan  
Munster Technological University



Succeeding Together

# Artificial intelligence (AI) Beyond the Hype



<https://nativevideotube.blogspot.com/>

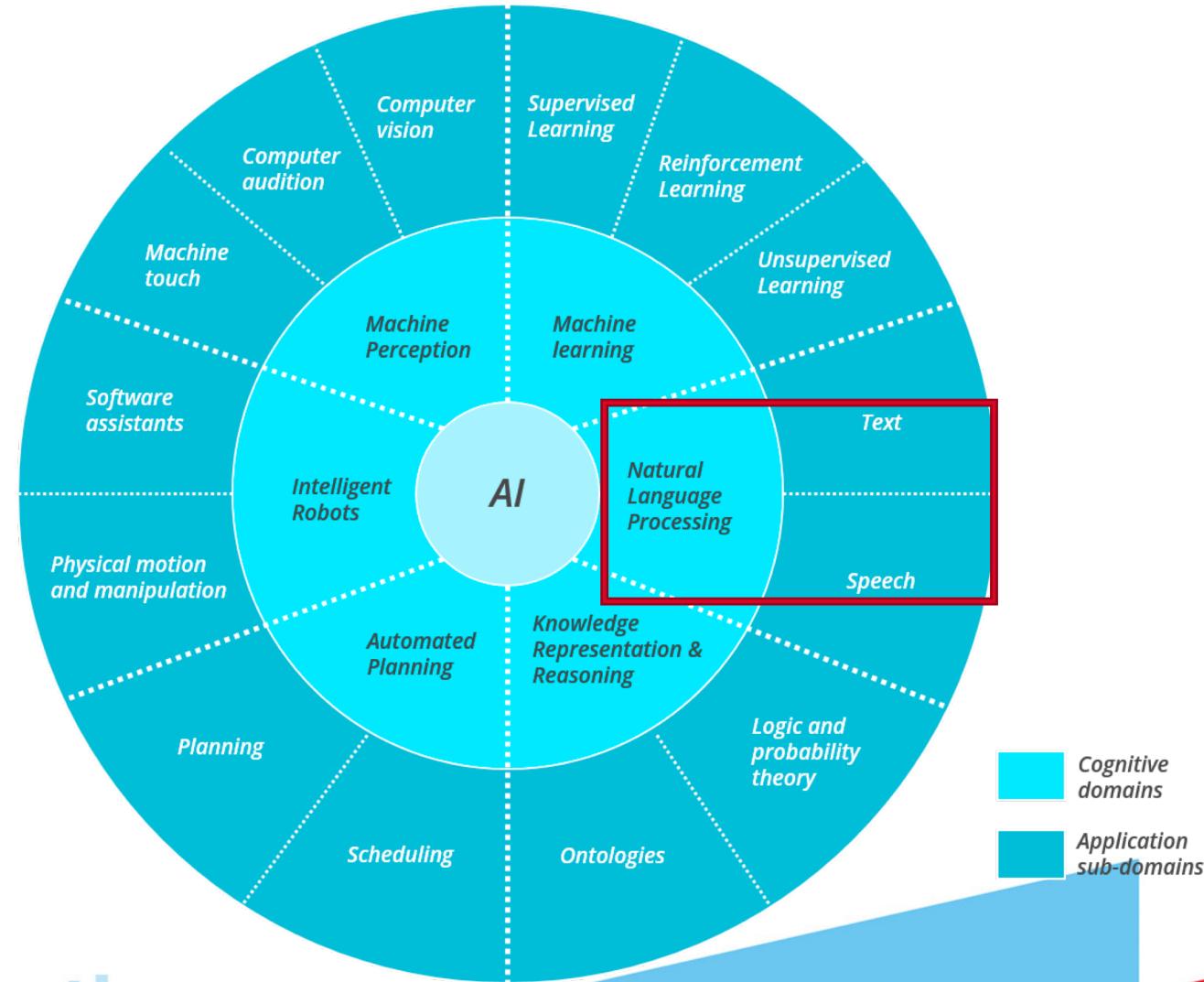
Graph from Tobias Bohnhoff

## Succeeding Together

# NLP - the language industry



**MTU**  
Ollscoil Teicneolaíochta na Mumhan  
Munster Technological University



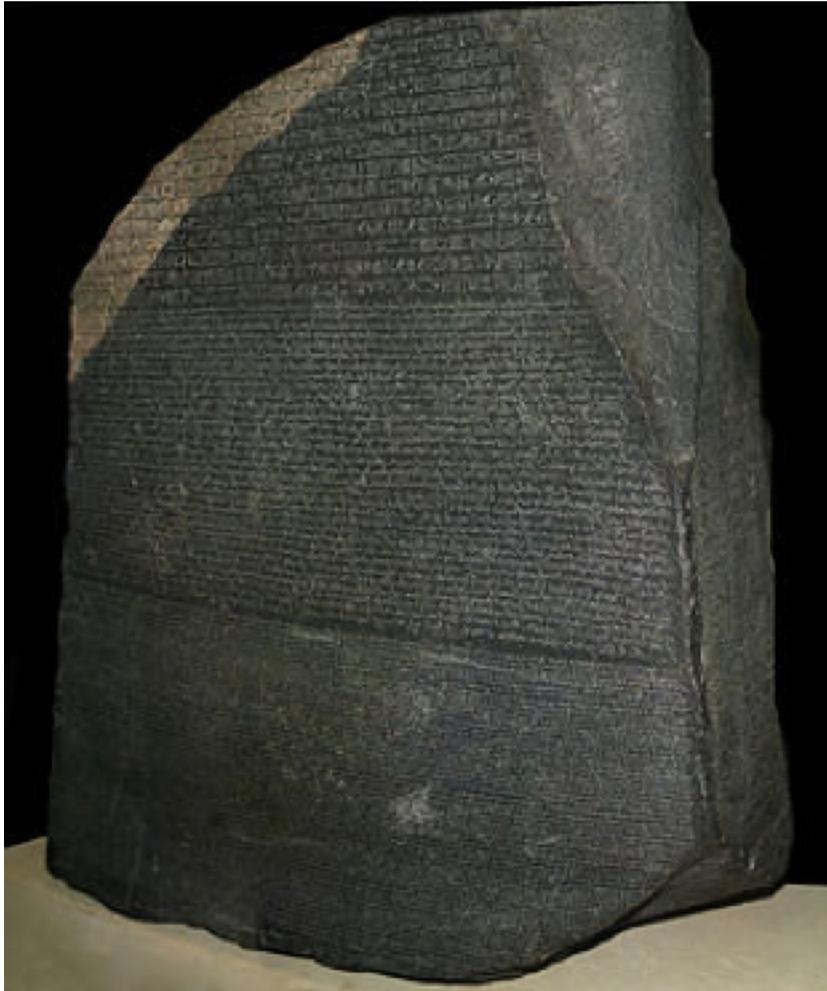
# Succeeding Together

# If you think the language industry is new



**MTU**

Ollscoil Teicneolaíochta na Mumhan  
Munster Technological University



For as far back as we can see,  
human has needed to  
communicate

-> Language industry paved  
the way for the renaissance of  
culture

# If you think the language industry is new, think again!



**MTU**

Ollscoil Teicneolaíochta na Mumhan  
Munster Technological University



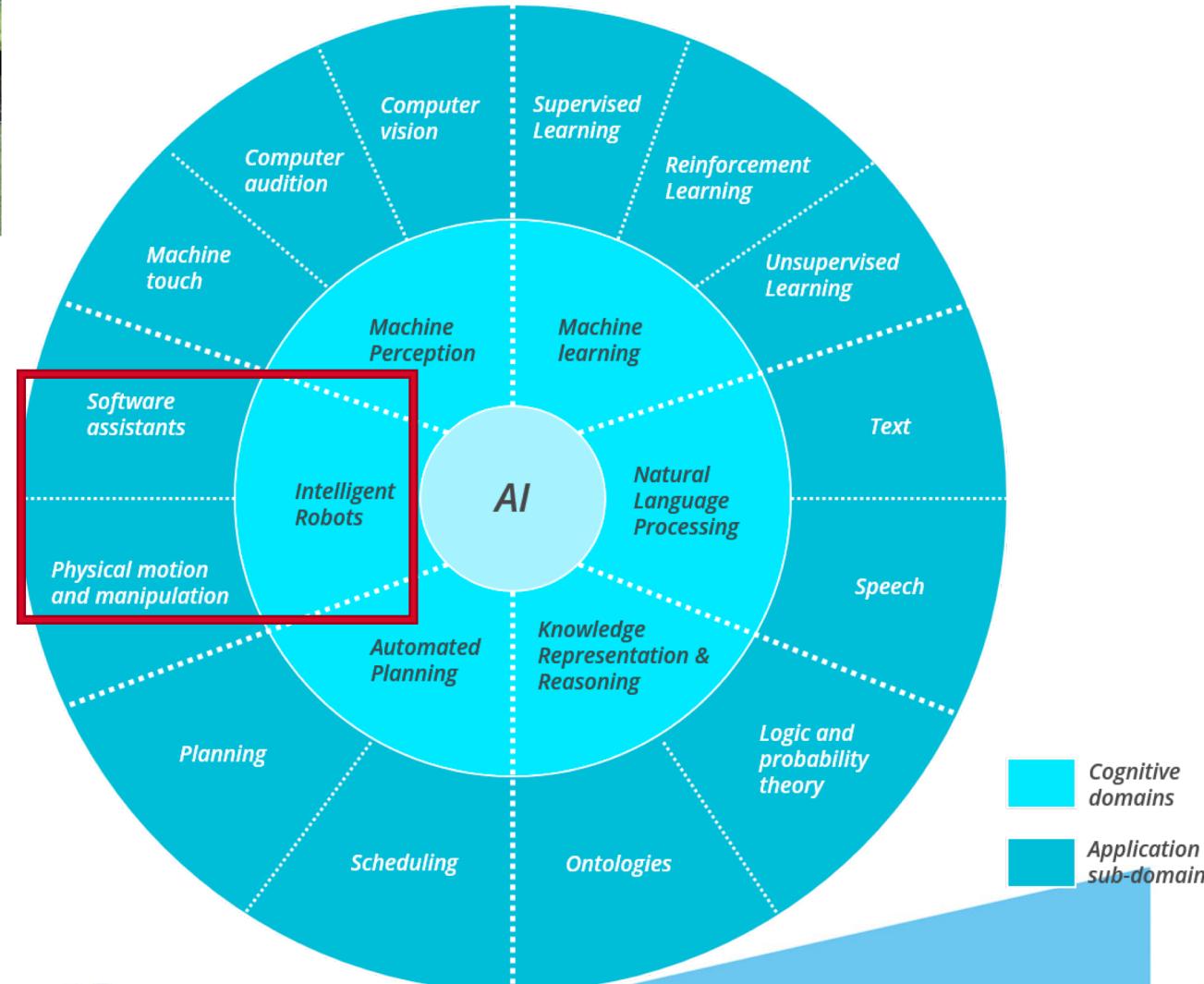
Succeeding Together

Rosetta Stone (British Museum)



**MTU**

Ollscoil Teicneolaíochta na Mumhan  
Munster Technological University



**Succeeding Together**

COMP9066 - Natural Language Processing - Haithem.afli@mtu.ie

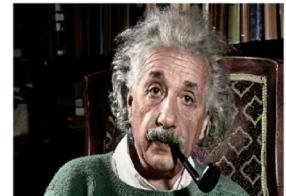
[www.mtu.ie](http://www.mtu.ie)

We're going to study what goes into getting computers to perform useful and interesting tasks involving human languages.

- **Text preprocessing** (stemmings, lemmatization, tokenization, stopword removal)
- **Probabilistic language models**
- **Text classifiers**
- **Word Embeddings**
- **Semantics**
- **Machine translation**
- **Dialogue Systems**



- **BERT** ([Google](#))
- **XLNet** ([Google/CMU](#))
- **RoBERTa** ([Facebook](#))
- **DistilBERT** ([HuggingFace](#))
- **CTRL** ([Salesforce](#))
- **GPT-2** ([OpenAI](#))
- **Megatron** ([NVIDIA](#))
- **ALBERT** ([Google](#))





**MTU**

Ollscoil Teicneolaíochta na Mumhan  
Munster Technological University

# COMP9066 - Natural Language Processing

## Module Workload

### Workload: Full Time

Workload Type	Workload Description	Hours	Frequency	Average Weekly Learner Workload
Lecture	Delivers the concepts and theories underpinning the learning outcomes.	2.0	Every Week	2.00
Lab	Application of learning to case studies and project work.	2.0	Every Week	2.00
Independent Learning	Student undertakes independent study. The student reads recommended papers and practices implementation.	3.0	Every Week	3.00
Total Hours				7.00
Total Weekly Learner Workload				7.00
Total Weekly Contact Hours				4.00

### Workload: Part Time

Workload Type	Workload Description	Hours	Frequency	Average Weekly Learner Workload
Lecture	Delivers the concepts and theories underpinning the learning outcomes.	2.0	Every Week	2.00
Lab	Application of learning to case studies and project work.	2.0	Every Week	2.00
Independent Learning	Student undertakes independent study. Student reads recommended papers and practices implementation.	3.0	Every Week	3.00
Total Hours				7.00
Total Weekly Learner Workload				7.00
Total Weekly Contact Hours				4.00

Succeeding Together

COMP9066 - Natural Language Processing - Haithem.afli@mtu.ie

[www.mtu.ie](http://www.mtu.ie)

# COMP9066 - Natural Language Processing



**MTU**

Ollscoil Teicneolaíochta na Mumhan  
Munster Technological University

Assessment Breakdown		%
Course Work		100.00%

Course Work				
Assessment Type	Assessment Description	Outcome addressed	% of total	Assessment Date
Project	Build a language model and use it in a given natural language processing application such as text generation. Produce a report that critically analyses the performance of the model.	1,2,3	50.0	Week 8
Project	Implement a machine model such as a neural model with vector-based representations for tasks of Machine Translation or Question answering. Assess the performance of the model using standard techniques such as BLEU or WER.	3,4	50.0	Week 12

No End of Module Formal Examination
-------------------------------------

Reassessment Requirement
<b>Coursework Only</b> <i>This module is reassessed solely on the basis of re-submitted coursework. There is no repeat written examination.</i>

The institute reserves the right to alter the nature and timings of assessment

Succeeding Together

COMP9066 - Natural Language Processing - Haithem.afli@mtu.ie

[www.mtu.ie](http://www.mtu.ie)

# Roadmap

- Introduction
  - Lecturer & module
  - **Topics**
  - Syllabus discussion
  - Books and Resources
- Lecture 1

# Natural Language Processing (NLP)

- Also known as
  - Computational Linguistics
  - Human Language Technology
- NLP is an interdisciplinary field
  - Computer science
  - Linguistics
  - Cognitive science, psychology, pedagogy, mathematics, etc.
- Applied natural language processing
  - Develop practical applications modeling human languages
- Theoretical computational linguistics
  - Focus on theoretical linguistics and cognitive science

# Natural Language Processing

- Applications

- Machine Translation (MT)
- Information Retrieval (IR)
- Automatic Speech Recognition (ASR)
- Optical Character Recognition (OCR)
- Automatic Summarization, Speech Synthesis, etc.

- Enabling Technologies

- Tokenization
- Part-of-Speech Tagging
- Syntactic Parsing
- Lemmatization
- Word Sense Disambiguation, etc.

# Natural Language Processing

- Rule-based/Symbolic Approaches
  - Linguists write rules that are applied by the machines
- Corpus-based/Statistical Approaches
  - Machines learn the “rules” from training data
    - Annotated data – supervised methods
      - Parallel Corpora: translated text collections
      - Treebanks: manually syntactically analyzed texts
      - Speech Corpora with transcripts
    - Unannotated data – unsupervised methods
    - Semi-supervised methods
  - Machine learning approaches are dominant in the field
- Hybrid Approaches
  - The best of **Smart**/Slow Humans and Dumb/**Fast** Machines

# Class Topics

- Basic text processing
- Finite state machines and word morphology modeling
- Language modeling
- Text Classification
- Sentiment analysis
- Conversational systems
- Machine translation
- A quick review of information retrieval.



**MTU**

Ollscoil Teicneolaíochta na Mumhan  
Munster Technological University

# Assignments

1. Basic text processing
2. Language Modelling
3. Text Classification
4. Machine translation / conversational systems

Succeeding Together

COMP9066 - Natural Language Processing - Haithem.afli@mtu.ie

[www.mtu.ie](http://www.mtu.ie)

# Roadmap

- Introduction
  - Lecturer & module
  - Topics
  - Syllabus discussion
  - **Books and Resources**
- Lecture 1

# Books and Resources

- **JM:** Daniel Jurafsky and James H. Martin. "Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition." (2nd or 3<sup>rd</sup> Edition).
- **MH:** Chris Manning and Hinrich Schütze, Foundations of Statistical Natural Language Processing, MIT Press. Cambridge, MA: May 1999.

<https://nlp.stanford.edu/fsnlp/>



**MTU**

Ollscoil Teicneolaíochta na Mumhan  
Munster Technological University

# Books and Resources

- Hands-On Machine Learning with Scikit-Learn and TensorFlow: Concepts, Tools, and Techniques to Build Intelligent Systems By Aurélien Géron.

Though not particularly dedicated to natural language processing, this practice-oriented book presents the most popular libraries that may be used for NLP and text analysis.

➤ Unix Lab

Succeeding Together



**MTU**

Ollscoil Teicneolaíochta na Mumhan  
Munster Technological University

# Programming Resources

We will be using Python 3 as our programming language in this module.

## Basic Python

[Python 3 Tutorial](#) – Clear and focused overview of Python 3 syntax, control structures, data structures etc.

[Video Python 3 Tutorials](#) – A set of very basic Python 3 video tutorials. More focused on beginners.

## NumPy and Pandas

DataCamp [NumPy Tutorial](#) – Accessible and easy to understand tutorial to get started with NumPy.

[NumPy Tutorial](#) – Short overview of NumPy and basic Python data structures. It also covers SciPy (which you don't need) and basic Matplotlib (which you will be covering later in the programme as part of visualization).

DataCamp [Pandas Tutorial](#) – Short and easy to understand tutorial on using Pandas.



**MTU**

Ollscoil Teicneolaíochta na Mumhan  
Munster Technological University

# Libraries and Tools

Spacy

NLTK

Transformers (Huggingface)

Gensim

Stanza

AllenNLP

Fast.ai

pattern

TextBlob

CoreNLP

CaMeL Tools ([https://github.com/CAMeL-Lab/camel\\_tools](https://github.com/CAMeL-Lab/camel_tools))

# The Jupyter Notebook



INSTALL PROJECT DOCUMENTATION BLOG DONATE



MTU

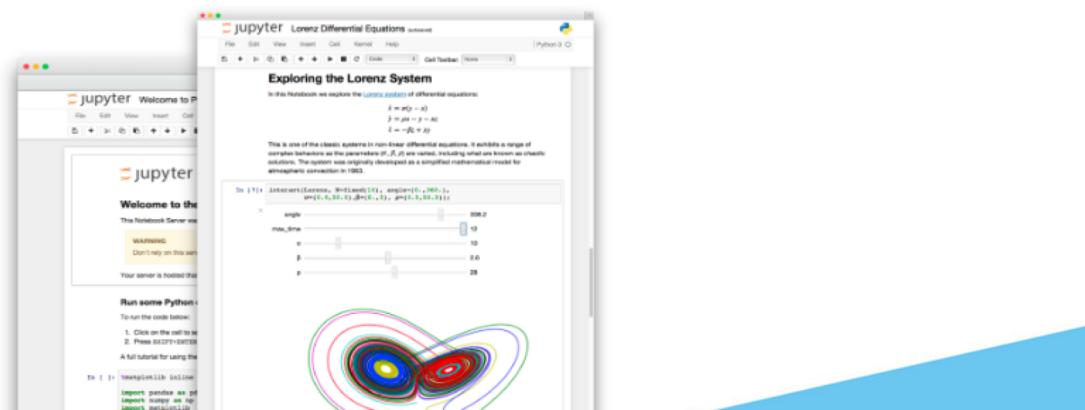
Ollscoil Teicneolaíochta na Mumhan  
Munster Technological University



<https://jupyter.org/>



The Jupyter Notebook is a web-based interactive computing platform that allows users to author data- and code-driven narratives that combine live code, equations, narrative text, visualizations, interactive dashboards and other media.



Succeeding Together

[www.mtu.ie](http://www.mtu.ie)



**MTU**

Ollscoil Teicneolaíochta na Mumhan  
Munster Technological University



# Discussion

Succeeding Together

[www.mtu.ie](http://www.mtu.ie)

# Roadmap

- Introduction
  - Lecturer & module
  - Topics
  - Syllabus discussion
  - Books and Resources
- **Lecture 1**

# Natural Language Processing



- We're going to study what goes into getting computers to perform useful and interesting tasks involving human languages.
- We are also concerned with the insights that such computational work gives us into human processing of language.

# Why Should You Care?

## Important trends

1. An enormous amount of knowledge is now available in machine readable form as natural language text
2. Conversational agents are becoming an important form of human-computer communication
3. Much of human-human communication is now mediated by computers

Very cool stuff! And with lots of commercial interest.

# Machine Translation

LIFESTYLE JULY 5, 2018 / 4:32 PM / 2 MONTHS AGO

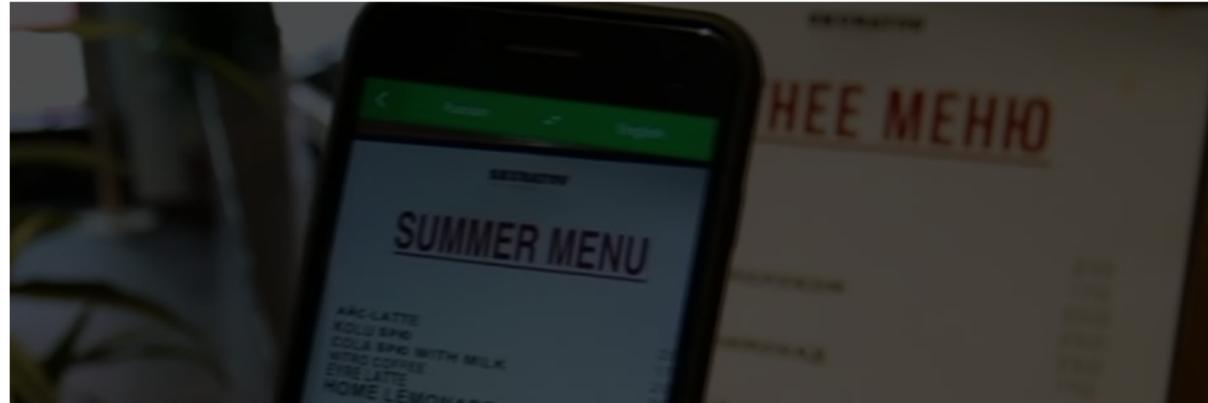


## At Russia's World Cup, Google Translate breaks language barriers

3 MIN READ



KAZAN, Russia (Reuters) - Soccer might be the most universal language on the planet. But when it comes to deciphering the Cyrillic alphabet or communicating with locals at the World Cup in Russia, the love of the game is sometimes not enough.



Succeeding Together

COMP9066 - Natural Language Processing - Haithem.afli@mtu.ie

[www.mtu.ie](http://www.mtu.ie)

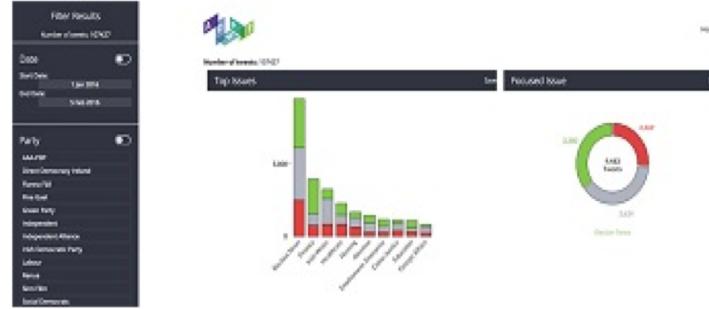
# Weblog Analytics

- Data-mining of Weblogs, discussion forums, message boards, user groups, and other forms of user generated media
  - Product marketing information
  - Political opinion tracking
  - Social network analysis
  - Buzz analysis (what's hot, what topics are people talking about right now).

# Social Media Analysis

ADAPT Centre for Digital Content Technology analyse Irish General Election activity

16 Feb 2016



Researchers at the ADAPT Centre for Digital Content Technology are working with RTE to monitor and analyse activity on Twitter related to the 2016 Irish General Election. The system that has been developed within ADAPT can quickly track the volume of tweets and gauge reaction and sentiment from the tweets in relation to election topics, parties and candidates. The project presents a graphical analysis of tweets that can be filtered by party or candidate allowing the user quickly understand and gain insights into the conversations on election topics taking place on Twitter.

<http://sma.adaptcentre.ie/ge16/#!/>

Haithem Afli, Sorcha McGuire, and Andy Way. 2017. Sentiment translation for low resourced languages: Experiments on irish general election tweets. In 18th International Conference on Computational Linguistics and Intelligent Text Processing.



**MTU**

Ollscoil Teicneolaíochta na Mumhan  
Munster Technological University

Succeeding Together

COMP9066 - Natural Language Processing - Haithem.afli@mtu.ie

[www.mtu.ie](http://www.mtu.ie)

# Information Extraction & Sentiment Analysis



**MTU**

Ollscoil Teicneolaíochta na Mumhan  
Munster Technological University



Size and weight

Attributes:

zoom



affordability



size and weight



flash



ease of use



- ✓ • nice and compact to carry!
- ✓ • since the camera is small and light, I won't need to carry around those heavy, bulky professional cameras either!
- ✗ • the camera feels flimsy, is plastic and very light in weight you have to be very delicate in the handling of this camera

Succeeding Together

# Sentiment Analysis (Aspect-based)



**MTU**

Ollscoil Teicneolaíochta na Mumhan  
Munster Technological University

CARA

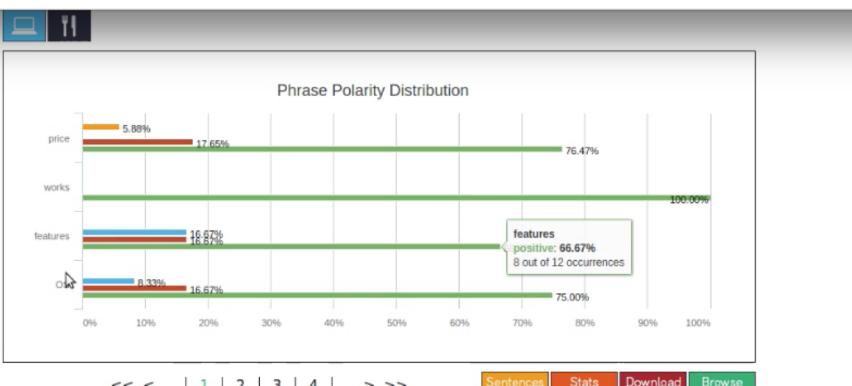
## CARA

Sentiment Analysis

The screenshot shows a list of four sentiment snippets:

- Boot time is super fast, around anywhere from 35 seconds to 1 minute.
- tech support would not fix the problem unless I bought your plan for \$150 plus.
- Positive Setup was easy.
- Did not enjoy the new Windows 8 and touchscreen functions.

Below the snippets are navigation buttons: << < | 1 | 2 | 3 | 4 | > >>. To the right are links for Sentences, Stats, Download, and Browse.



Succeeding Together

COMP9066 - Natural Language Processing - Haithem.afli@mtu.ie

[www.mtu.ie](http://www.mtu.ie)

The 2014 FIFA World Cup was the biggest event yet for Twitter with **672 million tweets**

Requested translation from Twitter (words)	 			 	Grand Total from all World Cup matches
	6,459,830	5,141,360		4,847,590	<b>85,047,110</b>

Top 3 languages

English

Portuguese

Spanish

- Source → Target traffic:
- EN → ES      13,614,450    (EN to all languages: 50,545,460)
- ES → EN      5,569,200    (ES to all languages: 10,609,420)
- PT → EN      1,831,750    (PT to all languages: 4,230,880)



**MTU**

Ollscoil Teicneolaíochta na Mumhan  
Munster Technological University

The 2014 FIFA World Cup was the biggest event yet for Twitter with **672 million tweets**

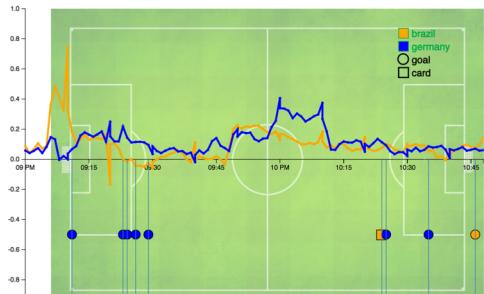
Requested translation from Twitter (words)		6,459	Grand Total from all World Cup matches	85,047,110
--	--	-------	--	------------

Top 3 languages

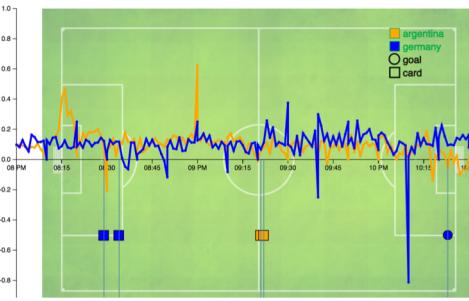
- Source → Target

- EN → ES 13,614,450 (EN to all languages: 50,545,460)
- ES → EN 5,569,200 (ES to all languages: 10,609,420)
- PT → EN 1,831,750 (PT to all languages: 4,230,880)

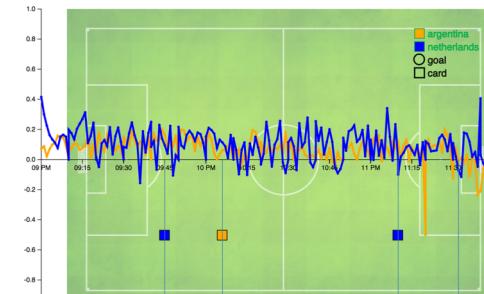
# UI: Sentiment pitch



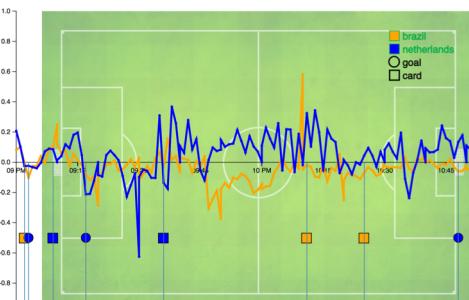
Semi-final: Germany 7-1 Brazil



Final: Germany 1-0 Argentina



Semi-final: Argentina 1-0 Netherlands

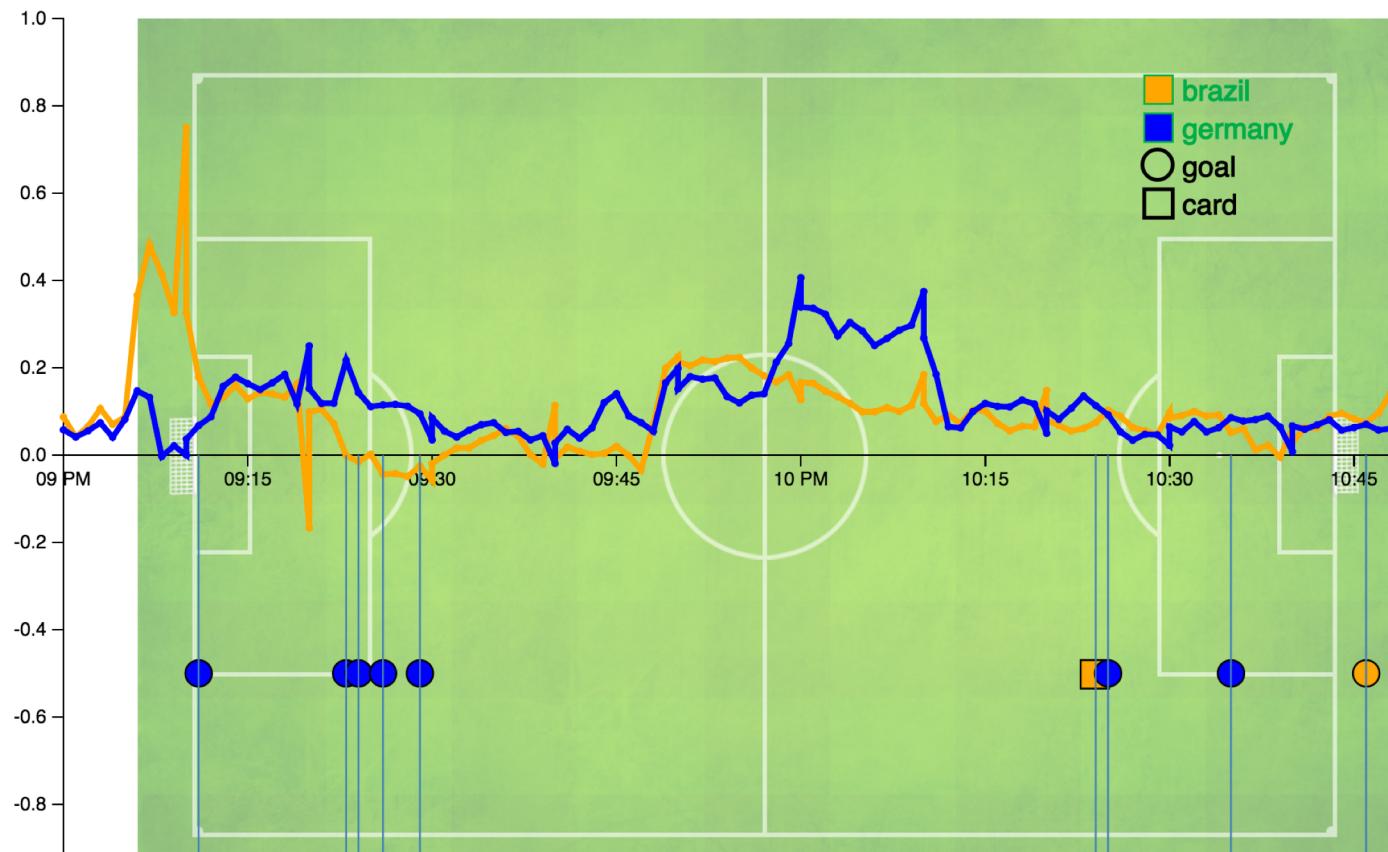


3<sup>rd</sup> Place: Netherlands 3-0 Brazil

Succeeding Together

[www.mtu.ie](http://www.mtu.ie)

# UGC Machine Translation - Braziliator



Semi-final: Germany 7-1 Brazil

Succeeding Together

[www.mtu.ie](http://www.mtu.ie)

# Applications

- First, what makes an application a *language processing application* (as opposed to any other piece of software)?
  - An application that requires the use of knowledge about human languages
    - Example: Is Unix wc (word count) an example of a language processing application?

# Applications

- Word count?
  - When it counts words: Yes
    - To count words you need to know what a word is.  
That's knowledge of language.
  - When it counts lines and bytes: No
    - Lines and bytes are computer artifacts, not linguistic entities

# Big Applications

- Question answering
- Conversational agents
- Summarization
- Machine translation



**MTU**

Ollscoil Teicneolaíochta na Mumhan  
Munster Technological University

# Big Applications

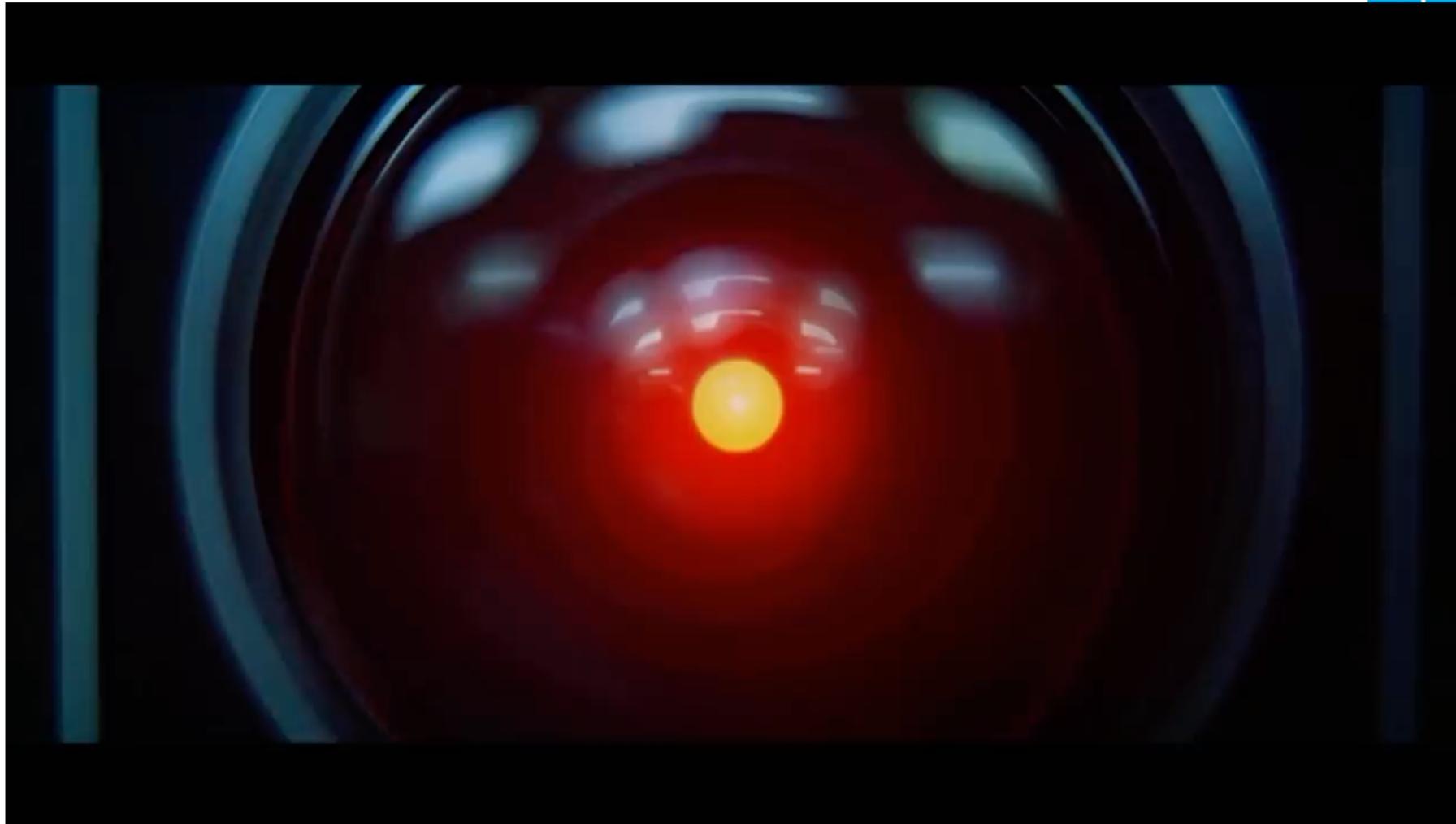
- These kinds of applications require a tremendous amount of knowledge of language.
- Consider the interaction with HAL the computer from **2001: A Space Odyssey**
  - Dave: *Open the pod bay doors, Hal.*
  - HAL: *I'm sorry Dave, I'm afraid I can't do that.*

# HAL 9000



**MTU**

Ollscoil Teicneolaíochta na Mumhan  
Munster Technological University



<https://www.youtube.com/watch?v=ARJ8cAGm6JE>

**Succeeding Together**

COMP9066 - Natural Language Processing - Haithem.afli@mtu.ie

[www.mtu.ie](http://www.mtu.ie)



**MTU**

Ollscoil Teicneolaíochta na Mumhan  
Munster Technological University

# What's needed?

- Speech recognition and synthesis
- Knowledge of the English words involved
  - What they mean
- How groups of words clump
  - What the clumps mean



**MTU**

Ollscoil Teicneolaíochta na Mumhan  
Munster Technological University

# What's needed?

- Dialog
  - It is polite to respond, even if you're planning to kill someone.
  - It is polite to pretend to want to be cooperative (I'm afraid, I can't...)

# Caveat

NLP has an **AI** aspect to it.

- We're often dealing with ill-defined problems
- We don't often come up with exact solutions/algorithms
- We can't let either of those facts get in the way of making progress



**MTU**

Ollscoil Teicneolaíochta na Mumhan  
Munster Technological University

# Course Material

- We'll be intermingling discussions of:
  - Linguistic topics
    - E.g. Morphology, syntax, discourse structure
  - Formal systems
    - E.g. Regular languages, context-free grammars
  - Applications
    - E.g. Machine translation, conversational systems



**MTU**

Ollscoil Teicneolaíochta na Mumhan  
Munster Technological University

# Topics: Linguistics

- Word-level processing
- Syntactic processing
- Lexical and compositional semantics
- Discourse processing
- Dialogue structure

# Topics: Applications

- Small
  - Spelling correction
  - Hyphenation
- Medium
  - Word-sense disambiguation
  - Named entity recognition
  - Information retrieval
- Large
  - Question answering
  - Conversational agents
  - Machine translation
- Stand-alone
- Enabling applications
- Funding/Business plans



**MTU**

Ollscoil Teicneolaíochta na Mumhan  
Munster Technological University

# Categories of Knowledge

- Phonology
- Morphology
- Syntax
- Semantics
- Discourse

- Each kind of knowledge has associated with it an encapsulated set of processes that make use of it.
- Interfaces are defined that allow the various levels to communicate.
- This usually leads to a pipeline architecture.

# Ambiguity

- Computational linguists are obsessed with ambiguity
- Ambiguity is a fundamental problem of computational linguistics
- Resolving ambiguity is a crucial goal



**MTU**

Ollscoil Teicneolaíochta na Mumhan  
Munster Technological University

# Ambiguity

- Find at least 5 meanings of this sentence:
  - I made her duck

# Ambiguity

- Find at least 5 meanings of this sentence:
  - I made her duck
- I cooked waterfowl for her benefit (to eat)
- I cooked waterfowl belonging to her
- I created the (plaster?) duck she owns
- I caused her to quickly lower her head or body
- I waved my magic wand and turned her into undifferentiated waterfowl



**MTU**

Ollscoil Teicneolaíochta na Mumhan  
Munster Technological University

# Ambiguity is Pervasive

- I caused her to quickly lower her head or body
  - **Lexical category:** “duck” can be a N or V
- I cooked waterfowl belonging to her.
  - **Lexical category:** “her” can be a possessive (“of her”) or dative (“for her”) pronoun
- I made the (plaster) duck statue she owns
  - **Lexical Semantics:** “make” can mean “create” or “cook”

# Ambiguity is Pervasive

- **Grammar:** Make can be:
  - **Transitive:** (verb has a noun direct object)
    - I cooked [waterfowl belonging to her]
  - **Ditransitive:** (verb has 2 noun objects)
    - I made [her] (into) [undifferentiated waterfowl]
  - **Action-transitive** (verb has a direct object and another verb)
    - I caused [her] [to move her body]



**MTU**

Ollscoil Teicneolaíochta na Mumhan  
Munster Technological University

# Ambiguity is Pervasive

- **Phonetics!**

- I mate or duck
- I'm eight or duck
- Eye maid; her duck
- Aye mate, her duck
- I maid her duck
- I'm aid her duck
- I mate her duck
- I'm ate her duck
- I'm ate or duck
- I mate or duck

# Dealing with Ambiguity

## Four possible approaches

1. Tightly coupled interaction among processing levels; knowledge from other levels can help decide at ambiguous levels.
2. Pipeline processing that ignores ambiguity as it occurs and hopes that other levels can eliminate incorrect structures.
3. Probabilistic approaches based on making the most likely choices.
4. Don't do anything, maybe it won't matter.
  1. *We'll leave when the duck is ready to eat.*
  2. *The duck is ready to eat now.*
    - Does the “duck” ambiguity matter with respect to whether we can leave?

# Models and Algorithms

- By **models** we mean the formalisms that are used to capture the various kinds of **linguistic knowledge** we need.
- **Algorithms** are then used to manipulate the **knowledge representations** needed to tackle the task at hand.

# Models

- State machines
- Rule-based approaches
- Logical formalisms
- **Probabilistic models**

# Algorithms

- Many of the algorithms that we'll study will turn out to be **transducers**; algorithms that take one kind of structure as input and output another.
- Unfortunately, ambiguity makes this process difficult. This leads us to employ algorithms that are designed to handle ambiguity of various kinds

# Why else is natural language understanding difficult?



**MTU**

Ollscoil Teicneolaíochta na Mumhan  
Munster Technological University

## non-standard English

Great job @justinbieber! Were SOO PROUD of what youve accomplished! U taught us 2 #neversaynever & you yourself should never give up either ❤️

# Why else is natural language understanding difficult?



**MTU**

Ollscoil Teicneolaíochta na Mumhan  
Munster Technological University

## non-standard English

Great job @justinbieber! Were SOO PROUD of what youve accomplished! U taught us 2 #neversaynever & you yourself should never give up either ❤️

## segmentation issues

the New York-New Haven Railroad

the New York-New Haven Railroad

# Why else is natural language understanding difficult?



**MTU**

Ollscoil Teicneolaíochta na Mumhan  
Munster Technological University

## non-standard English

Great job @justinbieber! Were SOO PROUD of what youve accomplished! U taught us 2 #neversaynever & you yourself should never give up either ❤️

## segmentation issues

the New York-New Haven Railroad  
the New York-New Haven Railroad

## idioms

dark horse  
get cold feet  
lose face  
throw in the towel

# Why else is natural language understanding difficult?



**MTU**

Ollscoil Teicneolaíochta na Mumhan  
Munster Technological University

## non-standard English

Great job @justinbieber! Were SOO PROUD of what youve accomplished! U taught us 2 #neversaynever & you yourself should never give up either ❤️

## segmentation issues

the New York-New Haven Railroad  
the New York-New Haven Railroad

## idioms

dark horse  
get cold feet  
lose face  
throw in the towel

## neologisms

unfriend  
Retweet  
bromance

# Why else is natural language understanding difficult?



**MTU**

Ollscoil Teicneolaíochta na Mumhan  
Munster Technological University

## non-standard English

Great job @justinbieber! Were SOO PROUD of what youve accomplished! U taught us 2 #neversaynever & you yourself should never give up either ❤️

## segmentation issues

the New York-New Haven Railroad  
the New York-New Haven Railroad

## idioms

dark horse  
get cold feet  
lose face  
throw in the towel

## neologisms

unfriend  
Retweet  
bromance

## world knowledge

Mary and Sue are sisters.  
Mary and Sue are mothers.

# Why else is natural language understanding difficult?



**MTU**

Ollscoil Teicneolaíochta na Mumhan  
Munster Technological University

## non-standard English

Great job @justinbieber! Were SOO PROUD of what youve accomplished! U taught us 2 #neversaynever & you yourself should never give up either ❤️

## segmentation issues

the New York-New Haven Railroad  
the New York-New Haven Railroad

## idioms

dark horse  
get cold feet  
lose face  
throw in the towel

## neologisms

unfriend  
Retweet  
bromance

## world knowledge

Mary and Sue are sisters.  
Mary and Sue are mothers.

## tricky entity names

Where is *A Bug's Life* playing ...  
*Let It Be* sold millions...  
... a mutation on the *for* gene ...

# Why else is natural language understanding difficult?



**MTU**

Ollscoil Teicneolaíochta na Mumhan  
Munster Technological University

## non-standard English

Great job @justinbieber! Were SOO PROUD of what youve accomplished! U taught us 2 #neversaynever & you yourself should never give up either ❤️

## segmentation issues

the New York-New Haven Railroad  
the New York-New Haven Railroad

## idioms

dark horse  
get cold feet  
lose face  
throw in the towel

## neologisms

unfriend  
Retweet  
bromance

## world knowledge

Mary and Sue are sisters.  
Mary and Sue are mothers.

## tricky entity names

Where is *A Bug's Life* playing ...  
*Let It Be* sold millions...  
... a mutation on the *for* gene ...

# Making progress on this problem...

- The task is difficult! What tools do we need?
  - Knowledge about language
  - Knowledge about the world
  - A way to combine knowledge sources
- How we generally do this:
  - Probabilistic models built from language data
    - $P(\text{"maison"} \rightarrow \text{"house"})$  high
    - $P(\text{"L'avocat général"} \rightarrow \text{"the general avocado"})$  low
  - Luckily, rough text features can often do half the job.

# Language Technology

making good progress



# MTU

Ollscoil Teicneolaíochta na Mumhan  
Munster Technological University

mostly solved

## Spam detection

Let's go to Agra!



Buy V1AGRA ...



## Part-of-speech (POS) tagging

ADJ ADJ NOUN VERB ADV

Colorless green ideas sleep furiously.

## Named entity recognition (NER)

PERSON ORG LOC

Einstein met with UN officials in Princeton

## Sentiment analysis

Best roast chicken in San Francisco!



The waiter ignored us for 20 minutes.



## Coreference resolution

Carter told Mubarak he shouldn't run again.

## Word sense disambiguation (WSD)

I need new batteries for my **mouse**.



## Parsing



I can see Alcatraz from the window!

## Machine translation (MT)

第13届上海国际电影节开幕...



The 13<sup>th</sup> Shanghai International Film Festival...

## Information extraction (IE)

You're invited to our dinner party, Friday May 27 at 8:30



Party  
May 27  
add

still really hard

## Question answering (QA)

Q. How effective is ibuprofen in reducing fever in patients with acute febrile illness?

## Paraphrase

XYZ acquired ABC yesterday

ABC has been taken over by XYZ

## Summarization

The Dow Jones is up

The S&P500 jumped

Housing prices rose



Economy is good

## Dialog

Where is Citizen Kane playing in SF?



Castro Theatre at 7:30. Do you want a ticket?

[www.mtu.ie](http://www.mtu.ie)

Succeeding Together

# Real Success: IBM's Watson



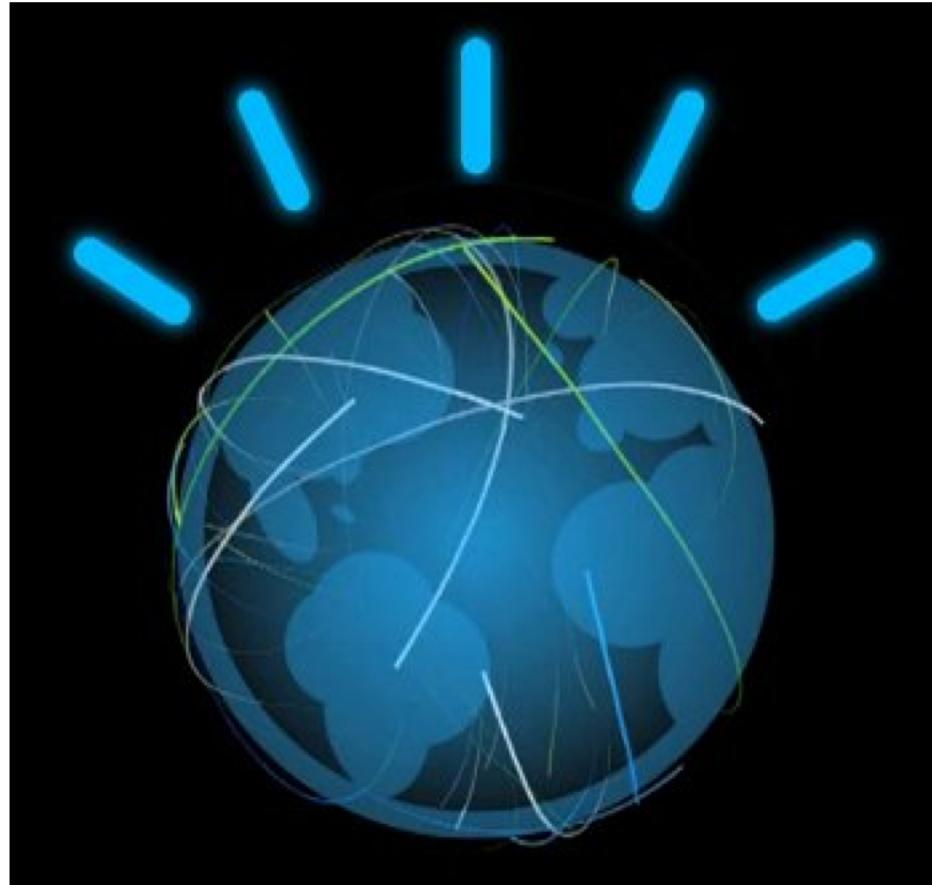
- Won Jeopardy on February 16, 2011!

WILLIAM WILKINSON'S  
“AN ACCOUNT OF THE PRINCIPALITIES OF  
WALLACHIA AND MOLDOVIA”  
INSPIRED THIS AUTHOR’S  
MOST FAMOUS NOVEL



Bram Stoker

# Real Success: Watson on Jeopardy



- [https://www.youtube.com/watch?v=WFR3lOm\\_xhE](https://www.youtube.com/watch?v=WFR3lOm_xhE)

- “There’s a big push to scale up machine learning to solve bigger and bigger problems, using more compute power and more data. As that happens, we have to be mindful of whether the benefits of these heavy-compute models are worth the cost of the impact on the environment.” **Dan Jurafsky**



- An off-the-shelf AI language-processing system can produce **1,400 pounds** of emissions – about the amount produced by flying one person **roundtrip between New York and San Francisco**.
- The **full suite of experiments** needed to build and train an AI language system from scratch can generate even more: **up to 78,000 pounds**, depending on the source of power. That's twice as much as **the average exhales over an entire lifetime**.

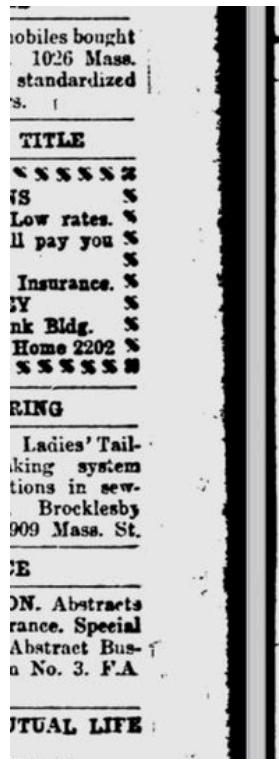
<https://hai.stanford.edu/blog/ais-carbon-footprint-problem>



- Adopting AI+NLP technologies is not an option

## For New-Age Businesses

- But Sustainable solutions will be needed.



**Horse vs. Automobile**

BEFORE you discard your horse and buy an auto it is well to think of the cost. Figure how much you spend for harness and then think of what new tires amount to. Figure up what it takes to feed Dobbin in a year and then think of gasoline, repairs and storage charges. Dobbin is worth what you paid for him two years ago, where's the man with an auto that can say the same? Come in and get a new harness instead of a new car and remember that Dobbin will take you through snow and mud as well as on good roads and that his carburetor is never out of order.

**Ed. Klein**  
732 Massachusetts Street

ly-company bond in the sum of fifty per cent (50%) of the contract price, for the faithful performance of the contract and specifications. He shall also give the State of Kansas a bond for the amount of the contract price, for the prompt payment for all materials and labor used on the work.

The County Commissioners reserve the right to reject any or all bids.

The successful bidder will be required to furnish his bonds and enter into a formal contract with the Commissioners, within ten (10) days from the date he is notified in writing that the contract has been awarded to him.

HERMAN BROEKER,  
County Clerk.

## Markets

Associated Press Market Report  
Kansas City, Mo., Aug. 3.—CATTLE—Receipts \$3,000; market steady to strong.

Prime Fed Steers \$9.50 to \$10.00.  
Dressed Beef Steers \$8.25 to \$9.50.  
Cows and Heifers \$4.50 to \$9.50.  
Suckers and Feeders \$6.25 to \$8.50.  
Bulls \$5.50 to \$7.00.  
Calves \$6.00 to \$10.25.  
HOGS — Receipts 7,000; market 10c to 15c higher.

this year will estimate that yield will be to an acre from twenty near Yellow near the summer as many as Park.

### HURT IN W

Mrs. Max Wil

Mrs. Max Wil wreck at Col but was not received a se and some bru of Lawrence in the accident on to Chicago

H. E. Wolf Methodist Ch to his home i



Farmers Hotel

Annapolis Royal, N.S.

Wm H Edwards Prope

1915 ad tries to convince people not to buy automobile and keep their horses

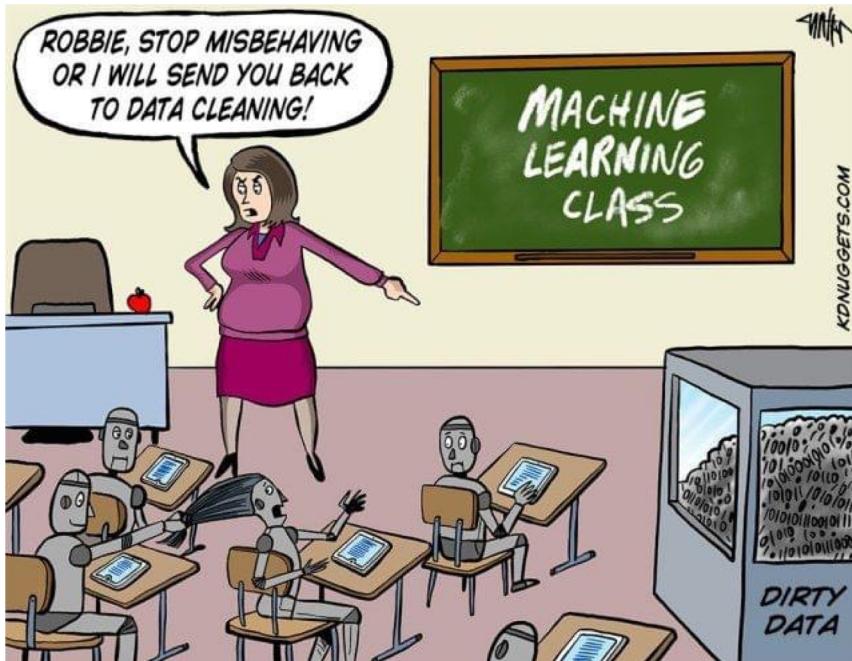


**MTU**

Ollscoil Teicneolaíochta na Mumhan  
Munster Technological University

# Next Week

- Basic text processing
  - Unix tools
  - Regular expressions





# Discussion

Some content was adapted from Speech and Language Processing - Jurafsky and Martin

# Thank You!

Haithem.afli@mtu.ie

Haithem Afli

@AfliHaithem

[www.mtu.ie](http://www.mtu.ie)