# NPTEL ONLINE CERTIFICATION COURSES

# DEEP LEARNING FOR NATURAL LANGUAGE PROCESSING

Lecture 01 : Introduction to the Course

**PROF. PAWAN GOYAL**

DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING
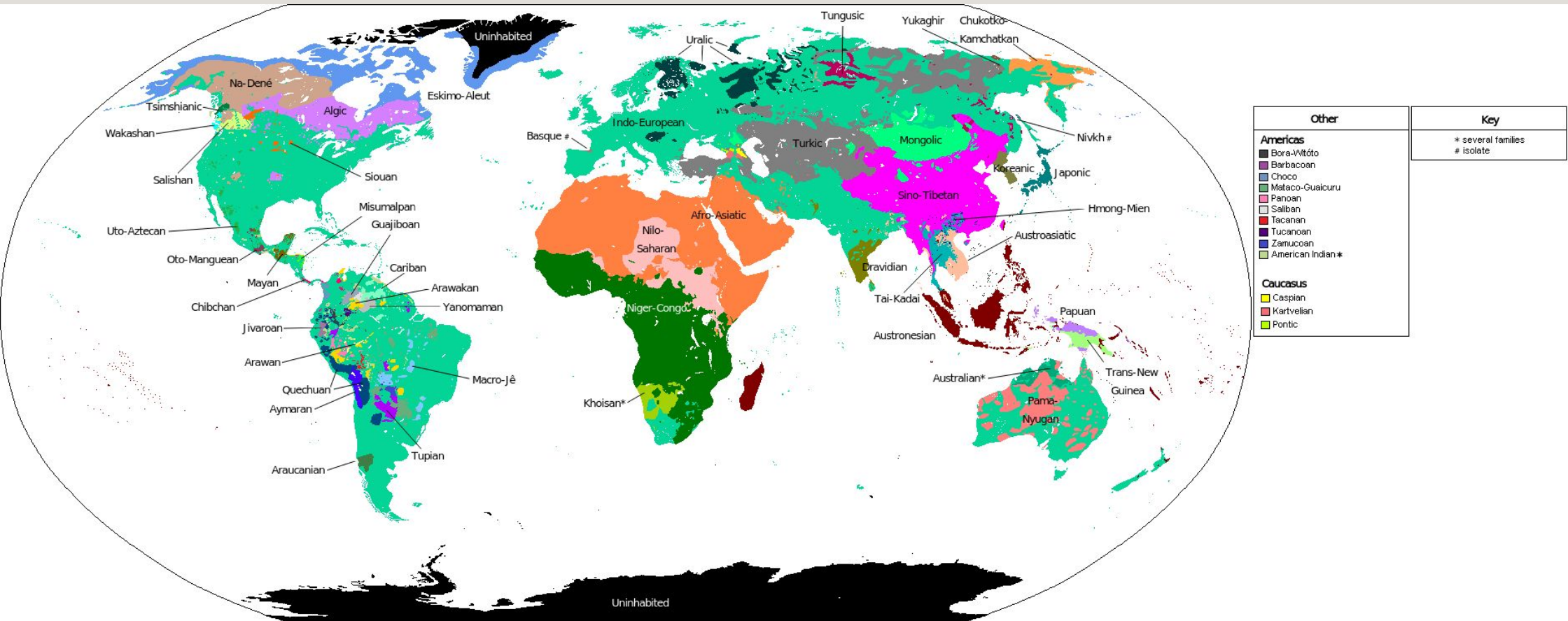
INDIAN INSTITUTE OF TECHNOLOGY KHARAGPUR

- Course Information

- What is NLP?

- Why Deep Learning for NLP?

- Course Content

# Course Information

- My Contact
  - Email: pawang@cse.iitkgp.ac.in
  - Webpage: http://cse.iitkgp.ac.in/~pawang/
  - Course Page: https://sites.google.com/view/dl4nlp-nptel/home

- Teaching Assistants (Inaugural Course)
  - Subhendu Khatuya
  - Pretam Ray

# Natural Language Processing



Natural Languages: Languages that evolved naturally through human use

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

# Natural Language Processing

## What is NLP?

- *Making computers understand what we write (or speak)*
- *Making computers write (and speak)*

The field of NLP attempts to design, implement and test systems

that process natural languages for practical applications
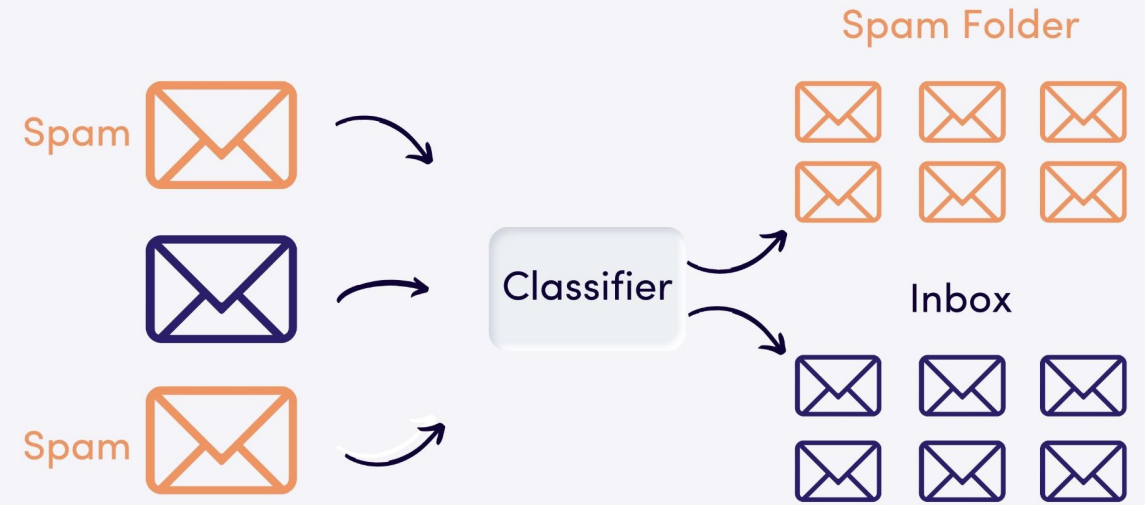
# NLP Applications: NLP is everywhere!

# NLP is everywhere!

# NLP is everywhere!

# NLP is everywhere!



What should i keep in mind while filing taxes? Tab

37/4000

What can i eat if i have a peanut allergy? Tab

30/4000

while creating an instagram post, what hashtags should i use? Tab

44/4000

# NLP is everywhere!

# Domain Specific Applications

# Why is NLP Hard? Language Ambiguity

**background** | ˈbakˌɡround |
noun

**1** [in singular] the area or scenery behind the main object of contemplation, especially when perceived as a framework for it: *the house stands against a background of sheltering trees.*
• the part of a picture or design that serves as a setting to the main figures or objects, or that appears furthest from the viewer: *the background shows a landscape of domes and minarets* | *the word is written in white on a red background.*
• a position or function that is not prominent or conspicuous: *after that evening, Athens remained* **in the background**.
• *Computing* used to describe tasks or processes running on a computer that do not need input from the user: *programs can be left running* **in the background**.
• *Physics* low-intensity radiation from radioisotopes present in the natural environment.
• unwanted signals, such as noise in the reception or recording of sound.

**2** the circumstances or situation prevailing at a particular time or underlying a particular event: *the political and economic background* | [as modifier] : *background information.*
• a person's education, experience, and social circumstances: *she has a background in nursing* | *a mix of students from many different backgrounds.*



Ok sir, next can you tell me a little about your background.

Sure, it's Mount Everest.

*Source: https://courses.cs.cornell.edu/courses/cs5740*

# Why is NLP Hard? Language Ambiguity



Source: https://courses.cs.cornell.edu/courses/cs5740

# Why is NLP Hard? Language Ambiguity

Let's try to decipher this weird conversation!

Rahul: *I saw a monkey with a banana.*

Computer: *That's gruesome!*

Rahul: *Why? What's so gruesome about seeing a monkey?*

Computer: *Oh I see! What else did you see with the banana?*

*In Natural Languages, ambiguity is the rule, not an exception*

*Example: Courtesy Dr. Monojit choudhury*

# NLP: Levels of Linguistic Structure

Discourse

Semantics

CommunicationEvent(e)        SpeakerContext(s)
Agent(e, Alice)                      TemporalBefore(e, s)
Recipient(e, Bob)

Syntax: Constituents

Syntax: Part of Speech

```
                              S
                             / \
                            /   \
                           /     VP
                          /     / \
                         NP    /   PP      .
                         |    /   /  \     |
                       Noun VerbPast Prep Noun Punct
```

Noun        VerbPast        Prep        Noun        Punct

Words

Alice   talked   to   Bob   .

Morphology

talk -ed   [VerbPast]

Characters

Alice talked to Bob.

Source: https://people.cs.umass.edu/~miyyer/cs685

# NLP Paradigms

We generally try to map problems to various (ML) paradigms

- Sentiment Analysis, news article groupings, etc. → Text Classification
- Named entity recognition, code-mixing, etc. → Sequence Labeling
- Machine Translation, summarization, chatbots, etc. → Text Generation

Timeline illustrating the progression of NLP from the 1950s

**Expert Systems and Statistical Models**

1. Rules and Ontology based Systems
2. Statistical Models
3. N-Grams commbined with Machine learning algorithms

**The Deep Learning Revolution**

1. Word2Vec, GLoVe, etc. word embeddings
2. Transfer Learning through pre-trained and fine-tuning
3. Attention Mecahnsim by Bahdanau et al.
4. Transformers by Vaswani et al.
5. BERT, GPT, and other models

| 1 | 1950s-1980s | 2 | 1980s-2000s | 3 | 2000s-2010s | 4 | 2010s-2020s | 5 | 2020s-now |

**Syntactic and Grammar-based**

1. Syntactic Structures by Noam Chomsky
2. ELIZA Chatbot
3. SHRDLU rule-based system

**Neural Models and Dense Representations**

1. Bengio et al.'s Dense Vector Representation
2. Mikolov et al.'s language Models based on Recurrent Networks
3. Pre-Trained Word Embeddings

**Era of LLLMs**

LLM

1. OpenAI releases GPT-2, GPT-3.5 and GPT 4
2. RLHF for alignment towards human values such as safety, groundedness, etc.
3. Open source LLMs and frameworks

Source: *Kamath, Uday, et al. "Large Language Models: A Deep Dive." (2024).*

# Why Deep Learning?



x = (0, ...., 0, 1, 0, ...., 0, 1, 0 ..... 0, 1, 0, ...., 0 , 1, 0 ,0, 1, 0, ...., 0, 0, 0 , ...., 0)

pw=the, w=dog, pt=NOUN, pt=DET, w=dog&pt=DET, w=dog&pw=the, w=chair&pt=DET

x = (0.26, 0.25, –0.39, –0.07, 0.13, –0.17) (–0.43, –0.37, –0.12, 0.13, –0.11, 0.34) (–0.04, 0.50, 0.04, 0.44)

**Word Embeddings**

| | |
|---|---|
| chair | (–0.37, –0.23, 0.33, 0.38, –0.02, –0.37) |
| on | (–0.21, –0.11, –0.10, 0.07, 0.37, 0.15) |
| dog | (0.26, 0.25, –0.39, –0.07, 0.13, –0.17) |
| | ... |
| | ... |
| the | (–0.43, –0.37, –0.12, 0.13, –0.11, 0.34) |
| | ... |
| | ... |
| mouth | (–0.32, 0.43, –0.14, 0.50, –0.13, –0.42) |
| | ... |
| | ... |
| gone | (0.06, –0.21, –0.38, –0.28, –0.16, –0.44) |
| | ... |

**POS Embeddings**

| | |
|---|---|
| NOUN | (0.16, 0.03, –0.17, –0.13) |
| VERB | (0.41, 0.08, 0.44, 0.02) |
| | ... |
| | ... |
| DET | (–0.04, 0.50, 0.04, 0.44) |
| ADJ | (–0.01, –0.35, –0.27, 0.20) |
| PREP | (–0.26, 0.28, –0.34, –0.02) |
| | ... |
| | ... |
| ADV | (0.02, –0.17, 0.46, –0.08) |
| | ... |

*Sparse vs. dense feature representations.* Two encodings of the information: current word is "dog;" previous word is "the" previous pos-tag is "DET."

Source: *Yoav Goldberg, Graeme Hirst. Neural Network Methods in Natural Language Processing, Morgan & Claypool Publishers (2017).*

These dense feature representations are used with various deep-learning architectures

Source: https://web.stanford.edu/~jurafsky/slp3.

# A timeline of the recent developments



Source: Alammar, J., & Grootendorst, M. (2024). Hands-On Large Language Models. O'Reilly.

# Change of NLP paradigms: Just use generation!



Sanh, Victor, et al. "Multitask Prompted Training Enables Zero-Shot Task Generalization." *ICLR 2022*

# Course Content (Weeks 1-6)

Background
- Introduction to NLP
- Introduction to Deep Learning and Representation Learning
- Word Representation: Word2Vec, Glove, FastText, Multilingual

Models and Architectures
- Recurrent Neural Networks: RNNs, LSTMs, Sequence to Sequence
- Attention Mechanism and Transformers: Attention in RNNs, Self-Attention in Transformers

Methods
- Pretraining: Self-supervised Learning objectives for Pretraining, ELMo, BERT, GPT, T5, BART, Fine-tuning

# Course Content (Weeks 7-12)

Tasks

- Question Answering, Text Summarization, Dialogs
- Domain and language-specific applications and challenges

Methods (LLMs)

- Towards building LLMs as chat assistants: Instruction Fine-tuning, Reinforcement learning from human feedback, Alignment techniques
- In-content learning, chain-of-thought prompting, Various LLMs
- Parameter Efficient Fine-tuning (PEFT), LoRA, QLoRA
- Handling Long Context, Retrieval Augmented Generation (RAG)

Conclusion

- Analysis and Interpretability, ethical considerations

## REFERENCES

Daniel Jurafsky and James H. Martin. 2024. Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition with Language Models, 3rd edition. Online manuscript released August 20, 2024. https://web.stanford.edu/~jurafsky/slp3.

Alammar, J., & Grootendorst, M. (2024). Hands-On Large Language Models. O'Reilly.

Yoav Goldberg, Graeme Hirst. Neural Network Methods in Natural Language Processing, Morgan & Claypool Publishers (2017).

Kamath, Uday, et al. "Large Language Models: A Deep Dive." (2024).

THANK YOU