# Natural Language Processing

Dr. Vaibhav Garg
Collegiate Assistant Professor

vaibhavg@vt.edu
Office# 479 (NVC)

Fall 2024

# Bio highlights and humble bragging

Ph.D. in Computer Science from NC State University

Passionate about solving problems having high social impact using NLP :)

Mentored both grad and undergrad students in research
        Sanjana Cheerla (Ph.D. student at NCSU) ; also an influencer LOL
        Ganning Xu (Undergrad at Georgia Tech)
        Arun Gaonkar (SDE 2 at Lexis Nexis)

Won Carla Savage Award for the most awesome Ph.D. student

Won best thesis award during Masters in CS

Recipient of NSA and NSF research grants

Accepted papers in top venues:
ACL and Communications of the ACM

Perhaps I'm the youngest :)

Don't like wearing formal but wa
to be taken seriously!!

# Bio highlights and humble bragging

I'm not a nerd and don't expect a bunch of nerds in my class!!

Enjoy taking breaks!!

Passionate bollywood dancer
Practice meditation
Badminton player
Excited about movies, filmmaking process, and songs

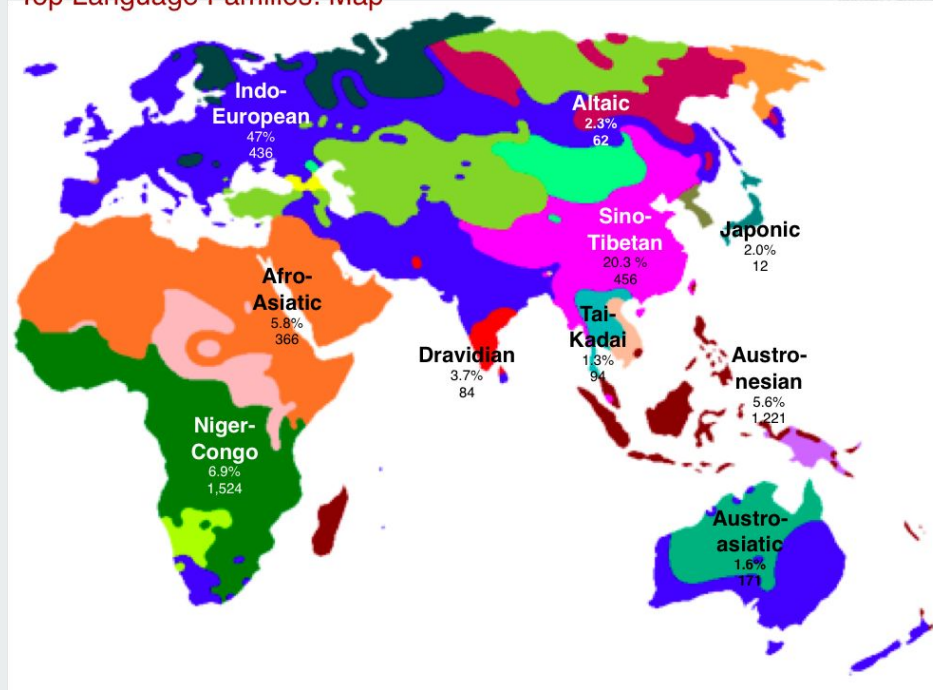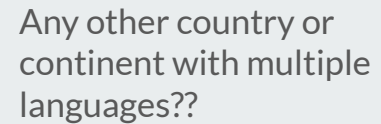Now your round of introduction??

# Diversity of languages



Top Language Families: Map

Indo-European 47% 436
Altaic 2.3% 62
Sino-Tibetan 20.3 % 456
Japonic 2.0% 12
Afro-Asiatic 5.8% 366
Tai-Kadai 1.3% 94
Dravidian 3.7% 84
Austro-nesian 5.6% 1,221
Niger-Congo 6.9% 1,524
Austro-asiatic 1.6% 171

https://triangulations.wordpress.com/2013/11/22/language-families/

# Variations across Indian languages



Any other country or continent with multiple languages??

https://x.com/indiainpixels/status/1299284220687167488

# What exactly is natural language?

Some naturally occurring phenomena?

A language created by someone? But then how is it natural?

How is it different from formal logic (if A implies B and B implies C; then A implies C)?

Let's discuss!

# What exactly is natural language?

*"A natural language is a human language, such as English or Standard Mandarin, as opposed to a constructed language, an artificial language, a machine language, or the language of formal logic. Also called ordinary language.*

*The theory of universal grammar proposes that all-natural languages have certain underlying rules that shape and limit the structure of the specific grammar for any given language."*

# Data data everywhere!!!



Evolving language–Phrases commonly used! GenZ slangs used over social media!

# NLP vs NLU vs NLG

"Alice is swimming against the **current**"     (current is noun)
"The **current** version of the paper is in the folder"   (current is adjective)

Is the above processing? or understanding? Or generation?

# NLP vs NLU vs NLG

"Alice is swimming against the **current**"
"The **current** version of the paper is in the folder"

Is the above processing? or understanding? Or generation?

Above is **understanding** based on grammar and semantics involved

NLG involves generating a new text (like QnA and chat bots)

NLP includes both – transforming freely flowing text into structured format which is easy to interpret

https://www.ibm.com/think/topics/nlp-vs-nlu-vs-nlg

# NLP



People think only deep learning and LLMs constitute NLP (**Completely wrong!!!!**)

# Course syllabus

Text normalization and preprocessing
N-gram language models
Naive Bayes, text classification, and sentiment analysis
Vector semantics and embeddings
Parts of speech and named entity recognition
Intro to transformers and LLMs

Fundamentals

Constituency and dependency parsing
Coreference
Semantic role labeling and information extraction
Lexicons for sentiment, affect, and connotation

Word senses and WordNet
Text summarization
Online argumentation

NLP Applications

# Smooth functioning of course

Will be using Canvas

Be active to check announcements

Don't spam with individual emails – first ask on Canvas

# Grading policy

Coding assignment 1 will be due in Sept            15%

Coding assignment 2 will be due in Oct             15%

Research paper presentation (Sept and Oct)    15%

Surprise quizzes (3 or 4) and seminars             15%

Final project                                                       40%

Yayyyy!! No exams :)

**Allowed:** Individual submissions late upto 2 days (total); Group submissions upto 1 day late (total)

# Academic integrity

Not allowed to copy the code or cheat in quizzes

Plagiarism check!!

Will also check for LLMs' redundant code

Severe reduction in grade!!!

# Lecture ethics

No usage of phones, tablets, and laptops

No talking and whispering (except in class discussions)

Feel free to interrupt and ask questions

Honesty in answers (discussions and quizzes)

Any questions related to the course structure and logistics?

# Challenges

Examples of each?



## Key Challenges in NLP

**Ambiguity**
Human language is inherently ambiguous, often relying on context and cultural nuances for accurate interpretation. Resolving this ambiguity remains a major challenge in NLP.

**Language Variability**
Languages exhibit variations across dialects, accents, and idiosyncrasies. Developing models that can handle such language variability is a complex undertaking.

**Sarcasm and Irony**
NLP struggles to capture the subtle nuances of sarcasm, irony, and other forms of figurative speech, which are prevalent in human communication.

**Lack of Contextual Understanding**
Understanding the context in which a word or phrase is used is crucial for accurate comprehension. NLP systems still face challenges in contextual understanding, leading to occasional misinterpretations.

# Ambiguity in NLP

**Lexical Ambiguity:** This type of ambiguity represents words that can have multiple assertions. For instance, in English, the word "back" can be a noun ( back stage), an adjective (back door) or an adverb (back away).

**Syntactic Ambiguity:** This type of ambiguity represents sentences that can be parsed in multiple syntactical forms. Take the following sentence: " I heard his cell phone ring in my office". The propositional phrase "in my office" can be parsed in a way that modifies the noun or on another way that modifies the verb.

**Semantic Ambiguity:** This type of ambiguity is typically related to the interpretation of sentence. For instance, the previous sentence used in the previous point can be interpreted as if I was physically present in the office or as if the cell phone was in the office.

**Metonymy:** Arguably, the most difficult type of ambiguity, metonymy deals with phrases in which the literal meaning is different from the figurative assertion. For instance, when we say "Samsung us screaming for new management", we don't really mean that the company is literally screaming (although you never know with Samsung these days ;) ).

# Multilingual



(a)      (b)

Code-mixed data

**Example I**

CODE-MIXED SENTENCE: is seat me girne ka koi chance nhi hai
ENGLISH TRANSLATION: there is no chance of falling down from this seat
REQUIRE CHANGES IN THE ENGLISH TRANSLATION?: No

**Example II**

CODE-MIXED SENTENCE: Thnks buds! Kabhi kabhi aajate hai achhe photos
ENGLISH TRANSLATION: Thank you buddy, sometime good photos are captured.
REQUIRE CHANGES IN THE ENGLISH TRANSLATION?: No

**Example III**

CODE-MIXED SENTENCE: Australia ke saath abhi jeete nahi hai, magar NZ ke saath final kaise jeetenge iss soch mein bhartiya yuvak on twitter.
ENGLISH TRANSLATION: Indian youth on twitter thinking that - We have not won against Australia yet, but how would we win final with NZ?
REQUIRE CHANGES IN THE ENGLISH TRANSLATION?: Yes

https://www.youtube.com/watch?v=lj0bFX9HXeE&t=3s

# Any takeaways?

# Acknowledgements