**Introduction to Large Language Models (LLMs)**
**Prof. Tanmoy Chakraborty, Prof. Soumen Chakraborti**
**Department of Computer Science & Engineering**
**Indian Institute of Technology, Delhi**
**Lecture 2**
**Introduction to Natural Language Processing**

Today we will try to introduce NLP within an hour, which is impossible. So what I will do, I will try to give you a very high-level overview of NLP. But as I mentioned earlier, you are strongly recommended to go through any NLP course.



Okay, so I will start off by showing this example. Okay, this is a recent news articles taken from Times of India on the assassination attempt on Donald Trump and this is the headline, right? Donald Trump's death. You may wonder possibly there are some mistakes in this title.

There are punctuation errors and so on and so forth. We will figure out later. But remember this and think whether the title is correct or not. Okay.

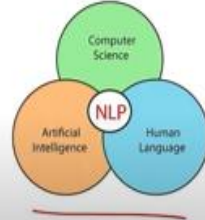Lec 02 | Introduction to Natural Language Processing

**Natural Language Processing**

What is a Natural Language?

Any language that has evolved naturally in humans through use and repetition without conscious planning or pre-meditation.

What is a Natural Language Processing?

A field of computer science, artificial intelligence and computational linguistics concerned with the interactions between computers and human (natural) languages.

But yeah, so let's go, let's proceed. What is NLP? So NLP is natural language processing. We all know what is a language. Language is a medium of communication. Natural language is essentially a language which has evolved without maintaining any grammar.

So it has evolved automatically, naturally and then the grammar has been introduced. So it is not like the artificial language like let's say Python or C++ or other types of programming languages where first the grammar was introduced and then the language was introduced. Natural language processing is a field of computer science, AI and computational linguistics which basically concerned with interactions between computers and human languages, right. So I mean this is of course controversial but generally we use this Venn diagram to identify the position of NLP within this space. We have CS, human language and AI and you know at the intersection of these three spaces we have this NLP, okay.

The Human Language

Home / India / More than 19,500 mother tongues spoken in India: Census

**More than 19,500 mother tongues spoken in India: Census**

There are 121 languages which are spoken by 10,000 or more people in India, which has a population of 121 crore, the report said.

https://indianexpress.com/article/india/more-than-19500-mother-tongues-spoken-in-india-census-5241056/

LCS | Tanmoy Chakraborty | LLMs: Introduction & Recent Advances

So why NLP is interesting because of the diversity language. Language has a lot of diversity. I mean, if you look at world's language, there are more than 6,000 languages across world. These are official languages by the way, right? So if you also consider dialects, nobody knows how many dialects are there in the world. But if you look at India, there are 19,000 languages, right? It's difficult to digest, but this is an official report.

It includes  the official languages as well as dialects. I am a Bengali, so I can tell you the amount of dialects that are there in West Bengal and in Bangladesh. If you go to West Bengal and Bangladesh, you will see crazy amount of variations. There are people from Bangladesh, they can tell you that if you go from east to west of Bangladesh, you see a huge variation. In fact, the western part of Bangladesh may not be able to understand the eastern part of Bangladesh.
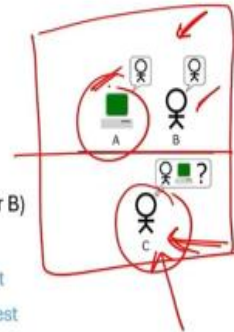
So lots of dialects and that essentially makes the entire business so difficult to process.

# Natural Language Processing

### Setup

- Two rooms, two humans, and a computer.
  - Room 1: One human C
  - Room 2: One computer (A) and one human (B)
- Response generated from room 2 (either by A or B)
- C has to figure out the source of the response
  - If C is successful → "A" failed the Turing test
  - Else, → "A" passed the Turing test

"Computing Machinery and Intelligence" proposed what is now called the Turing test.

Tanmoy Chakraborty — LLMs: Introduction & Recent Advances

Okay, so like computer science, NLP, I mean if you ask who is the father of NLP, of course one and only Turing, right, and we all know about Turing test. This is the famous paper computing machineries and intelligence. So in Turing test you essentially let's say you have a room. The room is partitioned and in one partition you have a this is human being and another partition you have a computer and human being.

So if there is a sound from this partition coming in, right, so the human being should, the task of the human being is basically to identify whether the sound is essentially generated by the machine or the human being, right. So if C is successful then A, the machine actually fails the Turing test. If C is not successful then A essentially passes the Turing test, okay.

## Natural Language Processing

In 1957, **Noam Chomsky**'s Syntactic Structures revolutionized Linguistics with '**universal grammar**', a rule-based system of syntactic structures

The father of modern linguistics | He is a laureate professor of linguistics at University of Arizona and an institute professor emeritus at MIT.

But of course when we talk about modern linguistics, computational linguistics, Noam Chomsky is considered as the father of modern linguistics, right. We all know the universal grammar proposed by Noam Chomsky, right.

We all know Chomsky's normal form, right. We will also talk about it maybe later if time permits. But he is a living legend.

## Why is NLP Challenging?

Ambiguity

Why is NLP challenging? This is challenging because of the ambiguity.



Let's look at this example again.

How many of you think that there is an error here? It should have been Donald Trump's death. What other things? Is this correct? Factually wrong? Yeah. So the Trump here is used as a verb. Trump card. So basically Donald Trump's death, okay, this was a very interesting example floated on twitter and my TA actually spotted this.

So I thought, this is a nice example to start with, okay and that's why NLP is really crazy.

# The Real Reason Why NLP is Hard

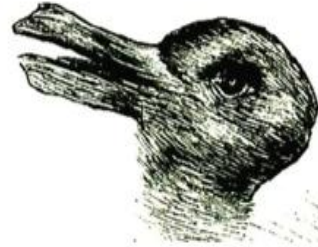Virat Kohli was on fire last night. He totally destroyed the other team.

There is another example right so if you look at this sentence. Virat Kohli was on fire last night, he totally destroyed the other team it doesn't mean that Virat Kohli was standing on fire, right so and we all know the meaning of this, right. So the literal meaning, the surface meaning of a sentence may not be same as the actual meaning of the sentence. And that is something where we need linguistics as well as computing to essentially teach the model to learn world's knowledge, to learn pragmatics, to learn discourse and so on and so forth.

## Ambiguity

Is ambiguity present in language only?

No, ambiguity is prevalent in every dimension!

Duck or Rabbit?

But it is not about language. where ambiguity is present, right? If you look at this picture, right? So somebody can say that this is a picture of a rabbit. Somebody can say that is a picture of a duck, right? So ambiguity is everywhere.



## Ambiguity in Language

- I saw a girl with a telescope.

  - I saw a girl with a bicycle.
  - I saw a bus with a telescope.

Who had the telescope?

OR

No ambiguity!

But we'll talk about ambiguity in language. Let's look at this sentence here.

Of course, you can say that these are very simple examples. I mean, human being can easily process these examples. But it's not as simple for a machine to process it. So look at this one. I saw a girl with a telescope.

So who had the telescope? The girl had the telescope or I had the telescope? It's not very clear. But of course, maybe if you think of it carefully, it seems that I had the telescope, but it is also possible that the girl had a telescope and I saw the girl. If you use the same syntax, the syntax that is present in this, the syntax is basically grammar. If you use the same syntax to write these two sentences, I saw a girl with a bicycle or I saw a bus with a telescope, things are very clear. There is no ambiguity because I can't see a girl through a bicycle.

So, it is obvious that the girl had the bicycle. Similarly, the bus did not have a telescope. Of course, bus is not a human being. So, it is very clear that I saw the bus through telescope. So the same syntax but the first case, the first sentence was a bit ambiguous.

If you do not know the semantics, I mean, it is very difficult to decipher the semantics whereas the next two sentences the semantics is very clear.



## Ambiguity in Language

- I saw a girl with a telescope.
- Mary had a little lamb.
    - Mary was physically bringing a lamb to a location, such as a farm or a home
    - Mary ate a lamb.

OR

Okay so the next example is this one. Mary had a little lamb. Okay so what do you think? Is this ambiguous? Yes. This is also ambiguous because in one case it is indicating that Mary had a lamb which was basically an animal, right? Whereas the other example, the other semantics can be that Mary basically ate a lamb.

Okay.

Lec 02 | Introduction to Natural Language Processing

## Ambiguity in Language

- I saw a girl with a telescope.
- Mary had a little lamb.
- Mujhe aapko mithai khilani padegi.

Let's look at some of the regional languages. Hindi, for example, Mujhe apko mithai khilani paregi. So, it is not very clear who is offering Mithai to whom. What do you think? Who is offering Mithai to whom? So this is called pragmatics. It's a very interesting aspect of NLP.

It has been a very interesting aspect of NLP before LLM hijacked the entire world of NLP. But pragmatics deals with the way you speak. The tone, the intonation, the hand movement, the face movement and so on and so forth. Based on that you determine  the meaning of this sentence, right?

## Ambiguity in Language

- I saw a girl with a telescope.
- Mary had a little lamb.
- Mujhe aapko mithai khilani padegi.
- I ate rice with spoon.
- I ate rice with curd.
- I ate rice with Rahul.

Similar surface structures but different interpretations!

Let's look at these three sentences now. I ate rice with spoon, I ate rice with card, I ate rice with Rahul, right? Same syntax, same syntax, right? Syntactically they are same, but semantically they are very different, okay? Spoon was used here as an instrument, right? Card was used as an auxiliary food with rice and Rahul here basically was used as somebody who was accompanying me.

So they are very different.

Now let's look at punctuation ambiguity. So if you look at the sentence "let's eat grandma" right versus "let's eat, comma grandma". If you look at the sentences again the only difference is in terms of this comma, right? But the meaning is completely different. So punctuation is very important i will show you some other examples, some other crazy examples look at this one. "A woman without her man is nothing", okay.

So we all understand the meaning of this, right? Now let's look at the beauty of punctuation. A woman without her man is nothing. This is a meaning that we can interpret easily. Now, if you add this colon here and a comma here, a woman colon without her man is nothing, right? So this is the power of punctuation.

What About This?

Is it a valid sentence?

Yes

Buffalo buffalo Buffalo buffalo buffalo buffalo Buffalo buffalo

The word *buffalo* has three senses:
1. Noun: Animal (plural is also buffalo)
2. Proper Noun: American State
3. Verb: To bully someone

Buffalo buffalo, whom other Buffalo *buffalo* buffalo, buffalo Buffalo buffalo

Dmitri Borgmann's *Beyond Language: Adventures in Word and Thought*. 1967.
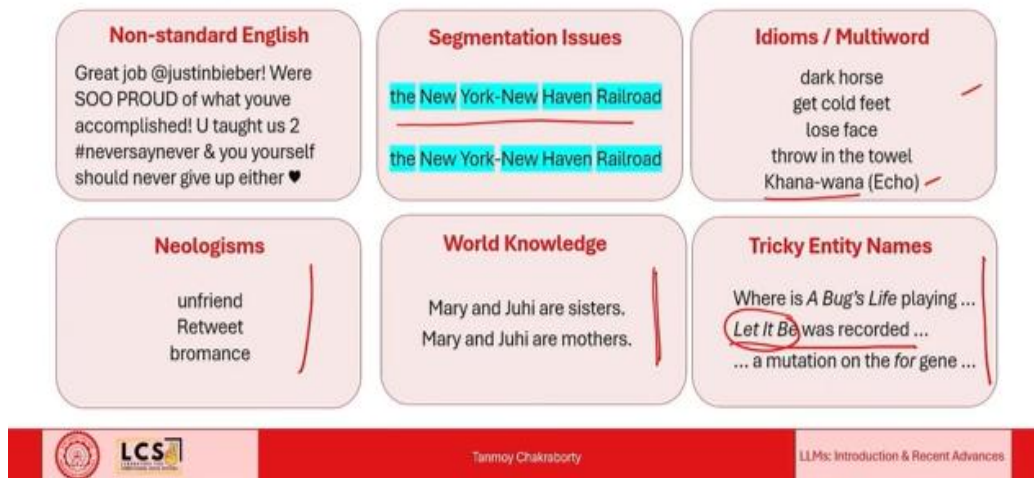
LCS | Tanmoy Chakraborty | LLMs: Introduction & Recent Advances

Now let's look at this one. This is my favorite example. I always start my NLP course with this example. So buffalo buffalo buffalo buffalo buffalo buffalo buffalo buffalo eight times the word buffalo is repeated. Do you think it's a valid sentence? This is a valid sentence and this is one of the crazy sentences that people always use to basically show that how anything is difficult. So if you look at this sentence, the word buffalo has three meanings.

Buffalo is used as an animal, buffalo is used as an American state and buffalo is used as a verb which means to bully someone. Clear, three meanings! Now look at it carefully, now look at this colors, color codes. Now look at this carefully, now I just added a few punctuation here. And there now look at it, buffalo buffalo meaning the "American states buffalo", right, whom other buffalo buffalo buffalo meaning So there are two sets of buffaloes here in American state. So one set of buffaloes are bullied by another set of buffaloes.

And they also bully them. So buffalo, buffalo, whom other buffalo, buffalo, buffalo? Buffalo, buffalo, buffalo. So this is a valid sentence.

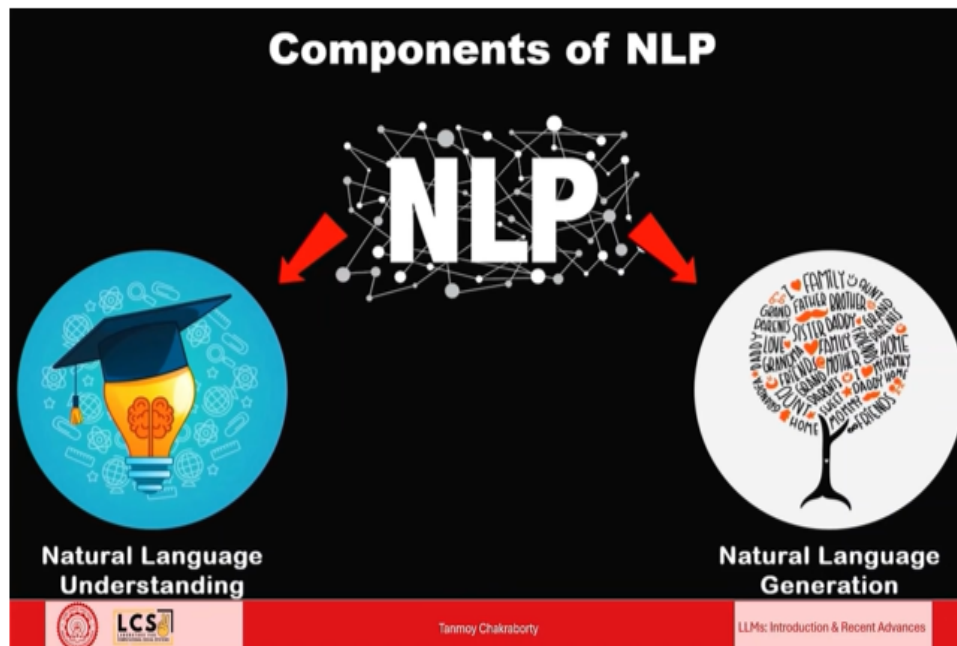# Why Else is Natural Language Understanding Difficult?

**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 / Multiword**

dark horse
get cold feet
lose face
throw in the towel
Khana-wana (Echo)

**Neologisms**

unfriend
Retweet
bromance

**World Knowledge**

Mary and Juhi are sisters.
Mary and Juhi are mothers.

**Tricky Entity Names**

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

Now why NLP is difficult? NLP is difficult because of specifically these issues, right? So we have non-standard English, right? And specifically in the era of social media, we see a lot of non-standard English like LOL, right? Retweet is also a grammatical, not a grammatical incorrect word, but it's not a dictionary word, okay? So non-standard English, then segmentation issue. If you feel that a space can be used to segment tokens, right? That may not be correct. For example, if you look at here, the New York, New Haven railway, right? If you tokenize it using space, you will get a different, you will get a separate, I mean, a very different and erroneous tokens, essentially.

The third problem is idioms. Multi word, okay. So multi word is a term that we use in NLP to indicate those terms which are whose surface meaning is very different from the actual meaning. For example hot dog okay. Surface meaning is a dog which is hot or who is hot whatever, right? But the actual meaning is basically a food. Similarly, all the phrasal verbs, right? Regional languages, if you look at phrase like "Khana-wana" and so on and so forth.

So it's basically those words which are not a part of dictionary, for example, unfriend, retweet, right? And so on and so forth. World's knowledge. World's knowledge is very important. For example, you know,  I think in 2014-2015 there was an article which essentially showed that one of the parsers, right, one of the parsers which was built based

on Stanford's corpus, it was not able to recognize Mitchell Obama as the wife of Barack Obama, right. It essentially returned Mitchell Obama as the wife of George Bush.

Okay, so, so was knowledge is very important. And it's, it's very difficult to incorporate was knowledge into into the model. And the next part, the tricky problem is the tricky named entity, right? Let's say, let it be was recorded. You need to understand that let it be is a single named entity, right? And It is referring to a song, right? But if these three words, let, eat, be appear in a different context, right? And then you consider let, eat, be as three different tokens, right? But if it appears in the context like this, then you have to identify or you have to mark let it be as a unique or single token, okay?



Okay, so in NLP, when you talk about NLP, NLP is categorized into two components or you know two aspects are there in NLP. One is natural language understanding, right, where we mostly try to understand a language, right, meaning let's say sentiment analysis, right, part of speech tagging, parsing and so on and so forth, right, whereas the other aspect is natural language generation. Right, where we talk about generating us language, generating a sentence, generating a phrase right chatbot basically comes under natural language generation where we generate the sentences tokens.

## NLP Layers

- Understanding the semantics is a non-trivial task.
- Needs to performs a series of incremental tasks to achieve this.
- NLP happens in layers.

| | |
|---|---|
| Pragmatics & Discourse | Study of semantics in context. |
| Semantics | Meaning of the sentence. |
| Parsing | Syntactic structure of the sentence. |
| Chunking | Grouping of meaningful phrases. |
| Part of speech tagging | Grammatical classes. |
| Morphology | Study of word structure. |

Increasing Complexity Of Processing

LCS  Tanmoy Chakraborty  LLMs: Introduction & Recent Advances

And this is the layer NLP layer that we follow generally. We start with morphology. Morphology deals with, I will talk about all these layers one by one but in a very very high level overview. Morphology is basically study of word structure, prefix, suffix, infix, circumfix and so on and so forth. Then we talk about  Now within syntax we have POS tagging, part of speech tagging.

We all know what is part of speech tagging. There is chunking. In chunking we essentially group different words into chunks. Noun chunk, verb chunk and so on. And then we have parsing. Parsing deals with the actual syntactic understanding of a sentence.

We try to build a parse tree, a grammar essentially from the sentence. then the most difficult part is semantics right. So morphology and syntax these two things essentially deal with the surface form of the sentence, surface form of the token, right? But when it comes to semantics, In the upper part of semantics, pragmatics discourse, these layers essentially connect NLP to the world, to the society. Because now we start thinking of pragmatics, context, world knowledge and so on. So semantics deals with understanding the meaning of it.

And then the last one is pragmatics and discourse. So pragmatics is all about context. Depending on the context, the meaning of the sentence can change. And discourse is basically about about a discussion right. Let's say there is a discussion happening between

two parties, multi parties right. We look at the entire context and then we respond accordingly okay.

So from bottom to top we basically move towards increasing complexity and the model will also become harder and harder.



Now this is called NLP Trinity. Okay, so you see a three-dimensional figure here, right? X-axis indicates different languages, right? We have Hindi, Marathi, English, French, etc. Y-axis indicates different tasks, right? Parsing, POS tagging, morphology, right? Semantical leveling and so on. So this is a task, this is a language, and this is an algorithm.

Okay. So CRF, HMM, Hidden Markov Model, Conditional Random Field, MEMM, right, Maximum Entropy Markov Model, there are other models, Deep Neural Network, LLMs, etc., right. So these are essentially models to solve these problems across these languages, okay.

## Word and Token

Word: Smallest sequence of phonemes of a spoken language that can be uttered in isolation.

Word Segmentation/Tokenization: Breaking a string of characters into a sequence of words.

Token: Smallest sequence of graphemes that are delimited with some predefined characters (space, comma, full-stop, etc.);

| | | |
|---|---|---|
| Ram, Shyam, and Mohan are playing. | ⇒ | [Ram] [,] [Shyam] [,] [and] [Mohan] [are] [playing] [.] |
| 21,53,010 COVID cases in India. | ⇒ | [21] [,] [53] [,] [010] [COVID] [cases] [in] [India] [.] |
| | | [21,53,010] [COVID] [cases] [in] [India] [.] ☑ |
| Check this out.(.https://www.abc.com | ⇒ | [Check] [this] [out] [.] [.] [.] [https] [:] [/] [/] [www] [.] [abc] [.] [com] |
| | | [Check] [this] [out] [...] [https://www.abc.com] ☑ |
| #GreatDayEver | ⇒ | [#] [Great] [Day] [Ever] |

Okay now let us look at this layer one by one.

I will again very very briefly talk about it. I will not talk about any algorithm to essentially solve this problem but if you, let us see if you go through Duravsky Martin's book you can easily understand it but also you can take any online lectures if needed. So the first and most important part is the word or a token, right? As I mentioned last class, I guess, a token can be a word, token can be a character, token can be a subword, right? And specifically in the context of LLM, tokenization is very, very important. Right, we will talk about different types of tokenization later on, white pair, what page, sentence page and so on and so forth. But tokenization is very important, don't assume that tokenization can be possible by you know using the space, you use the space and then you tokenize words, it's not that easy, okay. If you look at the first sentence here you may use punctuation like comma, full stop and space for tokenization.

But if you use the same way to tokenize this sentence it will produce an erroneous token because this entire number is a single token. Right? So this is important. Now, if you feel that, let's look at, I mean, if you consider full stop as a delimiter, for example, for tokenization, you fail here. Right? Because this URL is a single token. Right? You can't split the URL into multiple tokens based on this full stop.

Right? And similarly you know these days hashtags are very popular. So hashtag within a single hashtag there is no space at all, right? So great day ever if you want to tokenize it it would be very tough to tokenize.



## Morphology

- Field of linguistics that studies the internal structure of words
  - How they are formed
  - Their relationship to other words in the same language.
- It defines word formation rule from the root word.
- *Morpheme* is the smallest linguistic unit that has semantic meaning
  - *Example:*
    - *"Pre", "ed", "ing", "s", "es", etc.*
      - Dogs ⇒ dog + s (plural)
      - Going ⇒ go + ing (present participle)
      - Independently ⇒ independent + ly (Adverb)
        - ⇒ in + dependent + ly (Negation)
        - ⇒ in + depend + ent + ly (relying)
- pend: (verb) to remain undecided or unsettled.
        - ⇒ in + de + pend + ent + ly

**Morphology** is the study of words, how they are formed, and their relationship to other words in the same language. It analyzes the structure of words and parts of words, such as stems, root words, prefixes, and suffixes.

*Stemming*
*Lemmatization*

## Morphology

- English, Chinese, etc. are commonly referred as *morphologically-poor* language.
- Hindi, Turkish, Hungarian, etc. are termed as *morphologically-rich* language.

| English | Hindi | Linguistic property |
|---------|-------|---------------------|
| I will go. | मैं जाऊँगा। | |
| We will go. | हम जाएंगे। | Different morphological forms of word 'will go' in Hindi |
| You will go. | तुम जाओगे। | |
| He will go. | वह जाएगा। | |
| She will go. | वह जाएगी। | |

Okay so morphology. As I told earlier, morphology is basically a study of words, how they are formed, their relationship to other words in the same language. It also analyzes the structure of the word, parts of the word, suffix, stem, root words, prefix and so on.

Right, so if you have a morphological parser if you feed independently as an input to the parser. The parser will essentially return this independent plus ly. Independent can also be broken into in-dependent ly, right?

$$in + dependent + ly$$

So essentially we will see that parser may return this one in-de-pend-ent-ly.

$$in + de + pend + ent + ly$$

Right, so these are individual morphemes which are used to essentially form this particular token. How do you do this kind of tokenization? There are ways like I mean we can use finite state automata.

If you are aware of automata theory, use finite state automata to essentially and you have multiple rules and you use those rules for splitting a word into multiple morphemes. Now here in this context another aspect is very important. So essentially what morphology does? Morphology returns the stem of a word, the root word of a word, root word of a surface word. And it separates the other affixes like infix, prefix, postfix and so on. So in this context we need to understand the difference between stemming and lemmatization.

So stemming and lemmatization, these are two types of algorithms used for morphological analysis. Stemming returns a stemmed word or a root word which may not be a dictionary word. Whereas lemmatization returns a root word which will always be a dictionary word, right? Then the obvious question would be why do we need stemming, lemmatization can be useful I mean we can only use lemmatization why do we need stemming, the reason is that the reason behind using stemming is that stemming is very fast right, stemming is I mean compared to lemmatization because in order to do lemmatization you need to have a dictionary. Right, and dictionary is something which always changes over time right. So but stemming at least guarantees most cases, guarantees that if there are two different surface words if they are generated from the same lemma right stemming guaranteed.

Stemming always guarantees that their stemmed word is always the same. So you can always identify whether two surface words essentially came from the same root or not. So morphology, if you look at different languages, English, Chinese, etc., they are

morphologically poor languages. Right? Whereas languages like Hindi is morphologically rich.

Let's look at an example. So I will go, we will go, you will go, he will go, she will go. Right? Same kind of same surface sentences. Right? Now if you look at the Hindi translation of these sentences, ((25:34)) And so on and so forth. Right? So the complexity in Hindi language is that the verb has a gender. So in fact I have not been able to understand this properly but therefore languages like Hindi, Hungarian they are morphologically rich language.

## Syntax

**Syntax** concerns the way in which words can be combined together to form (grammatical) sentences.

| | | |
|---|---|---|
| **Pragmatics & Discourse** | *Study of semantics in context.* | |
| **Semantics** | *Meaning of the sentence.* | |
| Parsing | *Syntactic structure of the sentence.* | Increasing Complexity Of Processing |
| Chunking | *Grouping of meaningful phrases.* | |
| Part of speech tagging | *Grammatical classes.* | |
| **Morphology** | *Study of word structure.* | |

Syntax, so syntax concerns the way in which words can be combined to form a sentence. So POS tagging, chunking, parsing, these three steps basically are part of the syntactic analysis of a language.

# Parts-of-Speech (POS)

Grammatical class of the word.

| He | ate | an | apple | . |
|----|-----|----|-------|---|
| PRP | VBD | DT | NN | . |

**Tags**

PRP: Personal Pronoun
VBD: Verb, Past
DT: Determiner
NN: Noun, Singular, Mass
TO: *to*
IN: Preposition

- 45 tags in Penn Treebank tagset
- 146 tags in C7

**PoS disambiguation:**
o A word can belong to different grammatical classes.

| He | went | to | the | park | in | a | car | . |
|----|------|----|----|------|----|----|-----|---|
| PRP | VBD | TO | DT | *NN* | IN | DT | NN | . |

| They | went | to | *park* | the | car | in | the | shed | . |
|------|------|----|--------|-----|-----|----|----|------|---|
| PRP | VBD | TO | *VB* | DT | NN | IN | DT | NN | . |

POS tagging. Part of the tagging we all know.

How many part of speech you are aware of? Noun, pronoun. Adjective, verb, adverb, preposition, conjunction, interjection, article maybe. But in English, there are different types of tag sets. For example, if you follow Pantry Bank tag set, there are 45 English tags. Tags like PRP, personal pronoun, VRD, verb, past, DT, determiner, there are 45 such tags.

If you look at C7 tag set, there are 146 tags. So it's not as simple as we think of. But here also in POS tagging there is a problem of ambiguity. Let's look at these two sentences.

He went to the park in a car. They went to park the car in the set. So the word park here is used in the first sentence as a noun, in the second sentence as a verb. So if you think that if we can just create a dictionary or list where all possible words and the corresponding POS tags can be added and when we talk about POS, we just look at the dictionary and then fetch the relevant POS tag, that's not going to work. You need to understand the context again.

## Chunking

Identification of non-recursive phrases (noun, verb, etc.)

- He went to the Indian city Mumbai. ⇒
  [NP He] [VP went] [PP to] [NP the Indian city Mumbai]

- Mumbai green lights women icons on traffic signals earns global praise. ⇒
  [NP Mumbai green lights women icons] [PP on] [NP traffic signals] [VP earns] [NP global praise]

Chunking is the second step within syntax.

So chunking deals with essentially segregating different tokens into groups. For example here he, the word he is a noun phrase, went is a verb phrase, to is a prepositional phrase and the Indian city Mumbai is a noun phrase. Now let's look at this one. This is not that easy by the way. Let's look at this one.

Mumbai green lights women icons on traffic signals earns global praise. Right? So here Mumbai green signal women icons. This is a single noun phrase. On is a prepositional phrase, traffic signal is another phrase, another noun phrase and then on is a main part and then global phrase is essentially a noun phrase.

## Syntax Processing

Validate the grammatical structure of the sentence.

Let, vocabulary = [the, mango, he, eats, …]
He eats a mango. ⇒ ✅
He mango eats a. ⇒ ❌

- The sequence of words must follow the grammatical structure of the language to form a valid sentence.
  - Construct a parse tree.

Now once we are done with POS tagging and chunking, the next stage is parsing. Right and what is the purpose of a parsing? The purpose of a parsing is essentially given a sentence like this he eats a mango, right, your task would be to build this kind of tree, okay.

Let's assume that you already know the parser speech of all the tokens let's say he is PRP right eats is VBZ right A is determiner mango is NN right and then your task would be to group. This first level POS tags right into chunks, right. DT so determiner and noun they are combined to basically form the noun phrase, right. Then verb phrase and verb and noun phrase are combined to form the verb phrase and so on and so forth okay.

Syntactic Ambiguity

Now this is not simple let us look at this let us look at an ambiguous example, here, the same example I saw a girl with a telescope  Okay.

So if you look at the first level POS tagging, the first level POS tagging is the same, both the cases, right? But in this case, where I had the telescope, right? And via the telescope, I saw the girl. So what happened is that you know this noun phrase right and the preposition, they are combined together to form this PP prepositional phrase and PP, NP and verb. So verb noun phrase and prepositional phrase they are combined together to form the verb phrase right, whereas, here you see that this PP, NN, and determiner, they are combined together to form a noun phrase. So then the noun phrase and the verb, they are attached, and they form the verb phrase.

# Semantics

Semantics (and pragmatics) are the glue that connect language to the real world.

**Semantics** is concerned with the meaning of words and how to combine words into meaningful phrases and sentences.

- **Decompositional** – What the "components" of meaning "in" a word are

- **Ontological** – How the meaning of the word relates to the meanings of other words

- **Distributional** – What contexts the word is found in, relative to other words

In this way, we basically resolve the ambiguity problem. Now the next step is semantics and semantics as I mentioned earlier semantics is a step after which from semantics and after semantics we essentially try to connect language with the word okay. So semantics is concerned with meanings of words and how they are combined together to forms phrases and sentences. I mean if you tell me to take a class on semantics only, I can give a lecture on semantics on the entire class, I mean entire course, okay. It is very, very important, very, very kind of in-depth topic. But very high level semantics essentially can be classified or semantics of a sentence can be done in three ways.

The first one is decomposition, decomposition semantics, the second one is ontological semantics and the third one is distributional semantics.

## Decompositional Semantics

Decompositional Semantics Divides the Meaning of Words into Components

What are its strengths and weaknesses?

boy
$$\begin{bmatrix} +\text{human} \\ -\text{female} \\ -\text{adult} \end{bmatrix}$$

girl
$$\begin{bmatrix} +\text{human} \\ +\text{female} \\ -\text{adult} \end{bmatrix}$$

man
$$\begin{bmatrix} +\text{human} \\ -\text{female} \\ +\text{adult} \end{bmatrix}$$

woman
$$\begin{bmatrix} +\text{human} \\ +\text{female} \\ +\text{adult} \end{bmatrix}$$

Let us look at decomposition semantics. So in decomposition semantics, So what we do, we decompose a sequence into words and then we essentially try to decompose every word. Something like this. So here you see there are in this particular figure there is a boy, There is a girl, a man, and a woman.

Right. So we look at each and individual entity here separately and let's assume that these are the three features based on which we characterize each character or each entity, right. Human, female, adult. These are the three features based on which we characterize these four entities. So boy is a human, not a female, not an adult.

A girl is a human, is a female, right, not an adult. Man is a human, not a female, he is an adult, right. And woman is a female, human and an adult, right. So we look at individual entity and then we dig deeper into each of the entities and try to understand the semantics.

## Ontological Semantics

Ontological semantics says that the meaning of a word is its relationship to other words.

### The Basic (Ontological) Semantic Relations

- Synonymy—equivalence
  – <small, little>
- Antonymy—opposition
  – <small, large>
- Hyponymy—subset; is-a relation
  – <dog, mammal>
- Hypernymy—superset
  – <mammal, dog>
- Meronymy—part-of relation
  – <liver, body>
- Holonymy—has-a relation
  – <body, liver>

WordNet is a lexical resource that organizes words according to their semantic relations

- A graph
- A taxonomy
- An ontology

Ontological semantics. So ontology, some of you know what is an ontology, right? Can you give me an example of ontology? So WordNet is an ontology.

In WordNet, so WordNet is a resource or is a thesaurus developed through a consortium long back 1995, 1996 at that time, I think it was 6-7 years of long project, different versions of WordNet came up after that. So in WordNet what happens is that within a wordnet we have different tokens, right? Let's say let's say motor vehicle tractor right wagon a vehicle and so on and so forth golf cart and so on and so forth right, these are these are different entities. Entities are connected through different relations right. This is this is same as er diagram right database entity relationship diagram right. Here also there are different relations between entities is a right has part of right same as opposite to these are different relations right and these relations can be classified into these forms synonym, antonym, hyponym, myronym, holonym and so on and so forth right.

These are different types of relations and through these relations entities are connected right. For example, self-propelled vehicle is a wheeled vehicle, motor vehicle is a self-propelled vehicle, tractor is a self-propelled vehicle and so on and so forth. Car has part of car window and so on. So in ontological semantics then what we do given a sentence we look at the existence of individual tokens within the word net for example within the

ontology and then we return nearest neighbors of that token right. So for each individual token we have a set of neighbors which essentially characterize that token.

And then we look at the similarity between tokens in terms of the neighbors. If there is high overlapping between neighbors, then maybe two tokens are similar. But here the problem is that designing such an ontology is very tough. It requires a lot of manual labor and the way new words are coming in into the picture right it is not possible to come up with a complete ontology right where all possible entities, all possible words tokens are present okay.

## Distributional Semantics

The meanings of words can be derived from their distributional properties in large corpora of text. It relies on the context in which words appear.

**Example:** The meaning of the word "cat" can be inferred from the contexts it appears in, such as sentences where it co-occurs with words like "pet," "animal," "meow," and "feline."

The co-occurrence matrix

| | leash | walk | run | owner | pet | bark |
|------|-------|------|-----|-------|-----|------|
| dog | 3 | 5 | 2 | 5 | 4 | 2 |
| cat | 0 | 3 | 3 | 2 | 3 | 0 |
| lion | 0 | 3 | 2 | 0 | 1 | 0 |
| light | 0 | 0 | 0 | 0 | 0 | 0 |
| bark | 1 | 0 | 0 | 2 | 1 | 0 |
| car | 0 | 0 | 1 | 3 | 0 | 0 |

We essentially came up with this idea of distributional semantics, okay.

So this is the premise behind all recent models. For example, all recent word embedding methods like what to wear glove, we will discuss in the next week maybe, right. This is the premise behind this. So in distributional semantics, the idea is that a word is represented or the word is characterized by the surrounding words. right? The meanings of words can be derived from their distributional properties in large corpora of text, right? It relies on the context in which words appear. Let's look at an example, okay? Let's see you have words like this, dog, cat, lion, light, bark and so on, right? And you have a document, right? You have a corpus, from the corpus you build this matrix, and this matrix is called co-occurrence matrix.

So rows are unique tokens types, columns are unique tokens, it is a square matrix. Let us assume it is a square matrix and Of course this example is not a square matrix but ideally it's a square matrix okay and each entity within this matrix indicates number of times the word dog let's say this entity 5 indicates the number of times the word dog and the word walk co-appear together right What do you mean by co-appear together? how do you measure this we have a document and let's say you have a threshold, you have a sliding window of size 10, right. So you look at 10 words at a time. So you look at first 10 words and see whether they co-appear, I mean which words co-appear together. Then you move it by one step and then you look at again which words co-appear together and so on and so forth, right.

You determine the sliding window size and based on that you basically create this matrix. Now if you look at this matrix carefully, you will see the dog and cat, right? they co-appear with the word pit right they co-appear with the word owner right they co-appear with the word run right but they do not co-appear with the word bark. But still out of six words maybe in three to four words they co-appear together. That means they are semantically similar. But if you look at dog and car for example, they rarely co-appear except this owner. So simple way of looking at semantics in distribution semantics context is that each row acts as a vector for the token.

Now if you want to measure the similarity between two tokens, you essentially fetch two rows and you measure the similarity between two vectors using cosine similarity, dot product and you know other types of similarity measures. This is a very simple way, there are problems, we will discuss distributional semantics and the problems in the word embedding chapter.

## Pragmatics

Pragmatics considers [Thomas, 1995]:

- the negotiation of meaning between speaker and listener.
- the context of the utterance.
- the intention of the user.

- Context/World knowledge: An employee coming late to the office.
  - Utterance: Do you know what time is it?
  - Literal meaning: Are you aware of the current time? (Response: Yes, it is 12:30 PM)
  - Pragmatic meaning: Why are you coming so late? (Response: Reason for being late.)
- Intention:
  - Utterance: Can you pass the water bottle?
  - Literal meaning: Are you able to pass the water bottle? (Response: Yes, I can.)
  - Pragmatic meaning: Pass me the water bottle. (Response: Handover the water bottle)

Pragmatics deals with the context, pragmatics is essentially the negotiation of meaning between speaker and listener. So very nice phrase, I like this phrase a lot, this is taken from this paper, the negotiation of meaning between speaker and listener, I already give an example.

Right, do you know what time is it, right. The literal meaning is essentially are you aware of the current time and the answer should be the current time, whereas the pragmatic meaning is that you are essentially late, okay. If you look at this sentence, can you pass the water bottle, right? What would be the response? What would be the response? You just pass it, right? But if you say that yes I can, then stop, right? Then it is not a right answer.

It is a toughest, is a most difficult part of NLP, it essentially deals with dialogue, okay. Let's say, there is a sentence like this, a mother to John, go to school, it is open today, are you planning to bunk, father will be very angry, right. There is a dialogue between two parties and then you ask questions like what is open, bunk what, why the father is angry and so on and so forth.

So, you need to look at the entire discourse. context and then you respond, right. So the way we ask questions to ChatGPT, we have seen that ChatGPT has multiple sessions, right. So let's say within a session if you ask about NLP and suddenly you move to some political prompts, right, you will see the model fails terribly, right. So yesterday I was playing with ChatGPT and I was asking about an appropriate definition of NLP right but in a different context again the earlier context was different I was earlier I was talking about earlier I was discussing on some quantum theory etc and then suddenly I asked you know what is NLP? Then it returned that NLP is a neural language processing okay by the way if you search the word NLP on the internet right you see how many How many full phases are there for this acronym NLP? 13 to 14. Search with the word NLP and see how many full versions of NLP are available on the internet.

There are 13, 14 different versions of NLP. That is very crazy.

So let us look at very quickly, this is the last part of today's class. Let us look at some of the tasks that we are interested in and these tasks are generally used to evaluate the models, LLMs and other types of models. semantic role labelling. So semantic role labelling deals with this.

So let's say you have a sentence and you try to understand different roles of words. In this case, let's say John drove Mary from Delhi to Pune in his car. So the roles are agent, patient, source, destination. Now these roles are already given to you and your task would be to assign these roles to the words.

This is called semantic role labelling.

# Textual Entailment

Determine whether one natural language sentence entails (implies) another under an ordinary interpretation.

(*Ram hit Shyam with a hockey stick yesterday.* → *Shyam got hurt*)  ⇒ Positive TE .

(*Ram hit Shyam with a hockey stick yesterday.* → *Shyam did not get hurt*)  ⇒ Negative TE

(*Ram hit Shyam with a hockey stick yesterday.* → *Shyam got his first goal*)  ⇒ Non TE

The next one is textual entailment. entailment task right. You are given two sentences for example, let us say Ram hit Shyam with a hockey stick yesterday, Shyam got hurt okay. So whether the second sentence is entailed from the first sentence or not right. So in this case the first sentence entails second sentence right.

Therefore it is a positive entailment example. Whereas the second one, Ram hit Sam with a hockey stick yesterday. Sam did not get hurt. So negative entailment.

The last one, Sam got his first goal. So it has nothing to do with the first sentence. So this is non-entailment.

# Co-reference Resolution

- Two referring expressions used to refer to the same entity are said to co-refer.
- Determine which phrases in a document co-refer.

John shows Bob his Toyota yesterday. It's similar to the one I bought five years ago.

That was really nice, but he likes this *one* even better.

Co-reference resolution, it is actually very tough when the context becomes larger and larger. Let's say this one. John shows Bob his Toyota yesterday.

Look at this sentence, John shows Bob his Toyota yesterday. It is similar to the one I bought 5 years ago. That was really nice but he likes this one even better. Now here 'his' refers to whom? John or Bob. This is not clear. So in coreference resolution, we essentially resolve this problem. So there are different expressions like his, its, these, which it refers to is ambiguous and we essentially try to disambiguate it.

It, whether it is referring to Toyota or what. That was really nice. Now look at this that was. Here that refers to Toyota or to the one that I bought. five years ago. It is not very clear.

He refers to whom? John or Bob? Again it is not very clear. So in co-references resolution we resolve these problems. Information extraction we all know. Let us say you got an email. And you may have seen Google these days automatically block your calendar based on the email, right? They say you booked a hotel, right? You got a confirmation of your booking and you see that your calendar will be booked automatically by Gmail. This is information extraction.

## Information Extraction

Extraction of relevant piece of information.

- Named Entity Recognition (NER):
  - Identify names (Proper nouns)
    - [India]$_{Location}$ born [Sundar Pichai]$_{Person}$ is the CEO of [Google]$_{Organization}$ and its parent company [Alphabet]$_{Organization}$

- Relation Extraction:
  - Relation among entities
    - CEO(Sundar Pichai, Google), CEO(Sundar Pichai, Alphabet), Born-at(Sundar Pichai, India), ParentOrg(Alphabet, Google)

Information extraction also deals with you know identifying names named entities, right, relations between entities, right. So named entity in the recognition is again another task right. There are multiple types of named entity tags, person, location, organization, time, currency right and so on and so forth. There are multiple named entity tags and the task is essentially to identify named entities Relation extraction given let's say given two entities your task could be to extract relations. So these relations are either given or you can I mean your task can be to cure relations from the web right.

## Word Sense Disambiguation (WSD)

What does a word mean?

- The fisherman went to the *bank*.  ⇒  Financial bank or river bank?

  - The fisherman went to the *bank* to withdraw money.
  - The fisherman went to the *bank* to fish.

WSD is again very difficult task what sense is this disambiguation let's say for example the what bank whether bank refers to a financial sector Riverbank is something that we need to disambiguate.

Again Apple and so on.



Sentimental analysis all of you know right. Essentially the task is to categorize every sentence into positive sentiment, negative sentiment or neutral sentiment. Another related task is emotion detection. There are six different emotion tags, happy, sad, etc.

you categorize every utterance into emotions.

## Machine Translation

Given a sentence in the source language L1, convert it to the target language L2, such that the semantic (adequacy and fluency) is preserved.

Machine translation is very important task, it's one of the oldest task is NLP, right. I mean in this class I wanted to also talk about the history of NLP but I thought let's keep it. If you look at, if you go to Wikipedia and the history of NLP you see that machine learning, machine translation is a task we started early 1960. 1965 okay. And in 1965 you know that time I think English to Russian, Russian to English task was there and there was 65, 70 odd sentences and based on some rules they were able to successfully translate Russian to English and English to Russian and they thought that it was solved right but till date this is an open problem, machine translation. Okay when you talk about machine translation there are bias within the model okay for example if you translate she is a doctor with Google translator you see this kind of translation who a doctor have right English to Hindi now you use the same Hindi to again translated to English you see then the Google translator will say that he is a doctor right because Of course, I think that I mean, today's Google Translator has been updated, but it has been taken long time before.

But anyway, so the there is bias. But now look at this one. Okay. First, first pay khana sakth mana hai. This is a translation, eating carpet strictly prohibited.

## Summarization

Given a document, summarize the semantics (extract relevant information) in shorter length text.

LCS    Tanmoy Chakraborty    LLMs: Introduction & Recent Advances

Okay, so the next task is summarization, text summarization, where given a long paragraph or document, your task is to summarize it into a short paragraph, right, even very short kind of TLDR summary, one-line summary, right. There are three different types of summarizations. One is extractive summarization, where you identify important sentences from the document and you just copy-paste those sentences right, and you form the summary whereas abstractive summarization you identify important sentences as well as the semantics and you rephrase them right to generate the summary.

The third type is aspect based summary where given an aspect okay your task would be to only summarize document with respect to that aspect. Let us say you have a set of reviews of a product. You are only interested in getting a summary of the reviews with respect to let us say battery charge of a mobile. So the summary should be aspect based. And of course there are other types of summaries like extreme summarization like TLDR, multimodal summarization and so on.

## Question Answering

Answer natural language questions based on information presented in the repository.

**Factoid Questions**
- Question: Who is the author of the book Wings of Fire?
- Answer: A. P. J. Abdul Kalam

**List Questions**
- Question: What are the islands in India?
- Answer: Andaman Island, Nicobar Island, Labyrinth Island, Barren Island

**Descriptive Questions**
- Question: What is Greenhouse effect?
- Answer: The analogy used to describe the ability of gases in the atmosphere to absorb heat from the earth's surface.

Again an unsolved task right people earlier people have been thinking That this may be solved, but it's not a very easy problem. There are three types of questions that we are interested in factored questions right like who is the president, who is the father, who is the mother and so on and so forth right. Least type questions what are the islands in India right you basically produce a least and a descriptive question which is the difficult one. Fractured questions and least questions are easy to evaluate, but descriptive questions are difficult to evaluate.



## Dialog System and Chatbot

| | |
|---|---|
| C₁: | ...I need to travel in May. |
| A₁: | And, what day in May did you want to travel? |
| C₂: | OK uh I need to be there for a meeting that's from the 12th to the 15th. |
| A₂: | And you're flying into what city? |
| C₃: | Seattle. |
| A₃: | And what time would you like to leave Pittsburgh? |
| C₄: | Uh hmm I don't think there's many options for non-stop. |
| A₄: | Right. There's three non-stops today. |
| C₅: | What are they? |
| A₅: | The first one departs PGH at 10:00am arrives Seattle at 12:05 their time. The second flight departs PGH at 5:55pm, arrives Seattle at 8pm. And the last flight departs PGH at 8:15pm arrives Seattle at 10:28pm. |
| C₆: | OK I'll take the 5ish flight on the night before on the 11th. |
| A₆: | On the 11th? OK. Departing at 5:55pm arrives Seattle at 8pm, U.S. Air flight 115. |
| C₇: | OK. |

Dialogue systems you all know, right? Chatbot is an example of dialogue systems. Okay, so we will see how these tasks can be solved using large language model. Okay, thank you.