

Discourse Processing

Natural Language Discourse Processing

The most difficult problem of AI is to process the natural language by computers or in other words natural language processing is the most difficult problem of artificial intelligence.

If we talk about the major problems in NLP, then one of the major problems in NLP is discourse processing – building theories and models of how utterances stick together to form coherent discourse.

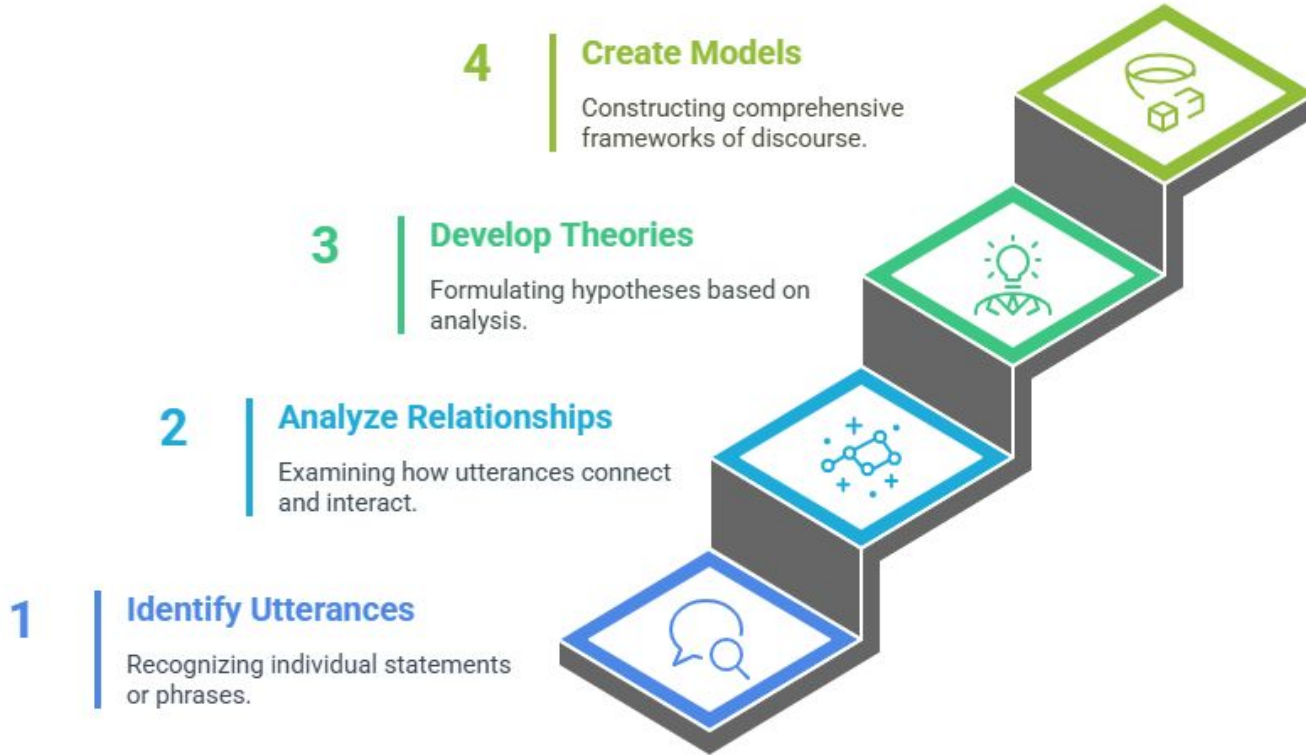
Actually, the language always consists of collocated, structured and coherent groups of sentences rather than isolated and unrelated sentences like movies.

These coherent groups of sentences are referred to as discourse.

Discourse in NLP is nothing but **coherent groups of sentences**. When we are dealing with Natural Language Processing, the provided language consists of structured, collective, and consistent groups of sentences, which are termed discourse in NLP.

The relationship between words makes the training of the NLP model quite easy and more predictable than the actual results.

Building Coherent Discourse Models



Coherent vs. Broken Discourse

Version A "The storm was approaching fast. Residents were advised to stay indoors. Emergency services were put on high alert to respond quickly."

Version B "Residents were put on high alert to respond quickly. The storm was stay indoors. Emergency services were approaching fast."

Rebuild the Story (Discourse Reconstruction)

1. She sent the results to her professor.
2. Finally, after hours of debugging, the code worked.
3. Ananya was working on her final project late at night.
4. She compiled the report early in the morning.

Identify Discourse Elements

“Ravi ordered a laptop from Amazon. It arrived in a damaged condition. He immediately contacted customer service. They promised to send a replacement.”

Tasks:

- What does “**it**” refer to?
- Who does “**he**” refer to?
- Who does “**they**” refer to?
- What's the topic shift, if any?

Identify Discourse Elements

“It” = the laptop

“He” = Ravi

“They” = Amazon customer service

The topic shift: From purchase ➡ damaged product ➡ complaint ➡ resolution.

1. Customer Service Chatbots (Case Study: Google's Meena, Meta's BlenderBot)

Problem: Chatbots often give generic or contextually irrelevant responses.

Discourse Challenge:

To maintain a meaningful conversation, a chatbot must track:

- **Who said what and when**
- **What the topic is**, even if it shifts subtly
- **User intentions** and emotional tone
- **Coreference**, e.g., understanding “he”, “it”, “they” from earlier dialogue

Example:

- User: "I ordered a phone last week, but it hasn't arrived."
- Chatbot: "I'm sorry to hear that. Can you share the order ID?"
- Later: User: "It was supposed to come yesterday."
 - ➔ The system must remember that **“it” refers to the phone**, and that this is a **delivery issue**, not a new order.

Fake News Detection (Case Study: DARPA's Semantic Forensics Program)

Problem: Automatically detecting fake news or misinformation.

Discourse Challenge:

Understanding coherence across a document is crucial. Fake news articles often:

- Lack logical progression
- Mix unrelated events
- Contain contradictions or omitted context

Example:

An article might claim:

"Vaccines have side effects. Also, a scientist in Norway discovered alien DNA in vaccines."

➡ NLP models must detect that these claims are **not logically connected**, possibly **fabricated**, and **lack coherence** — which involves discourse-level analysis.

Legal Document Summarization (Case Study: LexisNexis, CaseMine)

Problem: Automatically summarizing legal judgments.

Discourse Challenge:

Legal documents are structured with arguments, counterarguments, judgments, and precedents.

Understanding which sentences support, oppose, or conclude an argument is essential.

Example:

"The defendant argues X. However, the court noted Y. Therefore, the judgment is Z."

➡ A summary model must correctly associate "Z" as the final conclusion, and track the shifts in discourse from argument to judgment.

Discourse refers

- Discourse refers to any linguistic construction with multiple sentences.
- A disclosure is used in understanding and generating natural language.

A variety of text mining applications can be supported by discourse processing, which is a collection of Natural Language Processing (NLP) tasks used to extract linguistic structures from texts at different levels.

Identifying the conversational discourse's topic structure, coherence structure, coreference structure, and conversation structure is required for this.

Together, these structures can guide information extraction, sentiment analysis, machine translation, question answering, essay scoring, text summarization, and thread recovery.

Concept of Coherence

Coherence and discourse structure are interconnected in many ways.

Coherence, along with property of good text, is used to evaluate the output quality of natural language generation system.

The question that arises here is what does it mean for a text to be coherent? Suppose we collected one sentence from every page of the newspaper, then will it be a discourse? Of-course, not.

It is because these sentences do not exhibit coherence.

The coherent discourse must possess the following properties

The coherent discourse must possess the following properties

Coherence relation between utterances

The discourse would be coherent if it has meaningful connections between its utterances. This property is called coherence relation.

For example, some sort of explanation must be there to justify the connection between utterances.

Relationship between entities

Another property that makes a discourse coherent is that there must be a certain kind of relationship with the entities.

Such kind of coherence is called entity-based coherence.

✓ Example: News Article

Utterance 1: A massive fire broke out in the garment factory late last night.

Utterance 2: The fire department believes it was caused by a short circuit in the main panel.

✓ Example:

Utterance 1: The witness stated that the accused was present at the crime scene.

Utterance 2: The defense lawyer questioned the reliability of her statement.

Utterance 3: She had previously identified the wrong person in a similar case.

Discourse Structure

So far, we have discussed discourse and coherence, but we have not discussed the structure of the discourse in NLP. Let us now look at the structure that discourse in NLP must have. Now, the structure of the discourse depends on the type of segmentation applied to the discourse.

What is discourse segmentation ?

Well, when we determine the types of structures for a large discourse, we term its segmentation. The segmentation is a difficult thing to implement, but it is very necessary as discourse segmentation is used in fields like :

Information Retrieval,

Text summarization,

Information Extraction, etc.

Algorithms for Discourse Segmentation

We have different algorithms for Unsupervised Discourse Segmentation and Supervised Discourse Segmentation.

Unsupervised Discourse Segmentation

The class of unsupervised segmentation is also termed or represented as linear segmentation. Let us take an example to understand this discourse segmentation better.

Suppose we have a text with us, and the task is to segment the text into various units of multi-paragraphs. In the multi-paragraphs, a single unit is going to represent a passage of the text.

Now the algorithm will take the help of cohesion (that we have discussed above), and the algorithm will classify the dependent texts and tie them together using some linguistic devices. In simpler terms, unsupervised discourse segmentation means the classification and grouping up of similar texts with the help of coherent discourse in NLP.

The unsupervised discourse segmentation can also be performed with the help of lexicon cohesion. The lexicon cohesion indicates the relationship among similar units, for example, synonyms.

Supervised Discourse Segmentation

In the previous segmentation, there was no certain labeled segment boundary to separate the discourse segments.

But in the supervised discourse segmentation, we only deal with the training data set having a labeled boundary.

To differentiate or structure the discourse segments, we make use of cue words or discourse makers.

These cue words or discourse maker works to signal the discourse structure.

As there can be varied domains of discourse in NLP so, the cue words or discourse makers are domain specific.

Text Coherence

As we have previously discussed, the coherent discourse in NLP aims to find the coherence relation among the discourse text.

Now, to find the structure in discourse, we use lexical repetition, but by using this lexical repetition, we cannot satisfy the conditions of coherent discourse.

So, to prove such a kind of discourse relation, **Hebb** has proposed some solutions.

Suppose we have two kinds of related sentences, namely: S0 and S1.

Result

We can say that the second statement, i.e., S1 can be the cause of the first statement, i.e., S0. For example, **Rahul is late. He will be punished.**

In the above example, we can say that the first statement, S0, i.e., Rahul is late, has caused the second statement, i.e., S1, i.e., He will be punished.

Explanation

Similar to the result, We can say that the first statement, i.e., S0 can be the cause of the second statement, i.e., S1. For example, Rahul fought with his friend. He was drunk.

Parallel

By the term parallel, we mean that the assertion from the statement S_0 , i.e., $p(a_1, a_2, \dots)$, and the assertion from the statement S_1 , i.e. $p(b_1, b_2, \dots)$, the a_i and b_i is similar for all the values of i .

In simpler terms, it shows us that the sentences are parallel. For example, **He wants food. She wants money.** Both of the statements are parallel as there is a sense of want in both sentences.

Elaboration

Elaboration means that proposition P is inferring from both the assertions S_0 and S_1 . For example, **Rahul is from Delhi. Rohan is from Mumbai.**

Occasion

The occasion takes place when the change in the state is inferred from the first assertion S_0 , the final state is inferred from the statement S_1 , and vice-versa. Let us take an example to understand the relationship occasion better. For example,

Rahul took the money. he gave it to Rohan.

Building Hierarchical Discourse Structure

Let us now try to build a hierarchical discourse structure with the help of a group of statements. We generally create the hierarchical structure among the coherence relations to get the entire discourse in NLP.

Let us consider the following phrases and serially number them.

S1: Rahul went to the bank to deposit money.

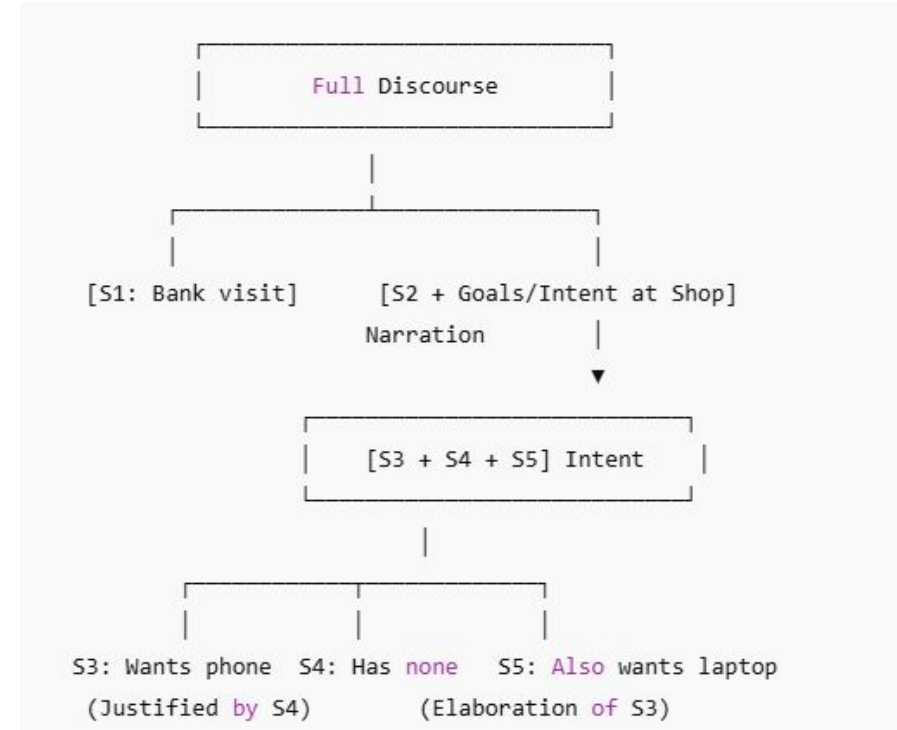
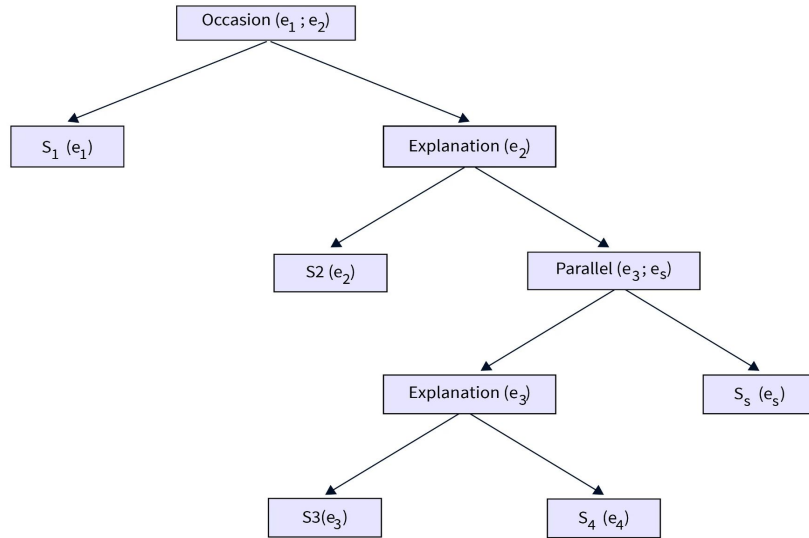
S2: He then went to Rohan's shop.

S3 : He wanted a phone.

S4 : He did not have a phone.

S5: He also wanted to buy a laptop from Rohan's shop.

Now the entire discourse can be represented using the below hierarchal discourse structure.



Reference Resolution

The extraction of the meaning or interpretation of the sentences of discourse is one of the most important tasks in natural language processing, and to do so, we first need to know what or who is the entity that we are talking about.

Reference resolution means understanding the type of entity that is being talked about.

By the term reference, we mean the linguistic expression that is used to denote an individual or an entity. For example, look at the below sentences.

- Rahul went to the farm.
- He cooked food.
- His farm was very big.

In the above sentences, Rahul, He, and His references. So, we can simply define the reference resolution as the task of determination of the entities that are being referred to by the linguistic expressions.

Terminology Used in Reference Resolution

Referring expression: The NLP expression that performs the reference is termed a referring expression. For example, the passage that we have talked about in the above section is an example of the referring expression.

Referent: Referent is the entity we have referred to. For example, in the above passage, Rahul is the referent.

Co-refer: As the name suggests, Co-refer is a term used for an entity if two or more expressions are referring to the same entity. For example, Rahul and He is used for the same entity, i.e., Rahul.

Antecedent: The term that has been licensed to use another term is termed antecedent. For example, in the above passage, Rahul is the antecedent of the reference He.

Anaphora & Anaphoric: The referring expression is termed anaphoric. Anaphora & Anaphoric can be said to be the term or reference used for an entity that has previously been introduced in the same sentence.

Discourse model: It is the model that has the overall representation of the entities that have been referred to in the discourse text. It also contains the relationship of the involved discourse in the NLP.

<https://www.scaler.com/topics/nlp/discourse-in-nlp/>

https://www.tutorialspoint.com/natural_language_processing/natural_language_discourse_processing.htm