

Pre - Print Version

# The Role of NLP in Coreference Resolution in Sindhi Text

Saira Baby Farooqui<sup>1\*</sup>, Noor Ahmed Shaikh,<sup>2</sup> Samina Rajper<sup>3</sup>

**Abstract-** Finding terms in a text that relates to the same thing is a significant difficulty in natural language processing (NLP). We call this procedure “coreference resolution.” This task is crucial for many NLP applications, such as information extraction, text summarization, and machine translation. Even though coreference resolution has been thoroughly explored in English and other commonly used languages, the difficulties presented by the Arabic language call for unique strategies catered to its unique linguistic and grammatical traits. Sindhi is a highly inflected language with a rich derivational and inflectional morphology system, flexible word order, and intricate agreement patterns. These linguistic features introduce complexities that impact traditional coreference resolution techniques. Additionally, Arabic exhibits variations across dialects, further complicating the task due to differences in syntactic structures and lexical choices.

## INTRODUCTION

This research presents an overview of coreference resolution in the Sindhi language. We delve into the linguistic phenomena that contribute to the challenges in this task, such as pronominalization, noun-pronoun relationships, and the use of definite and indefinite articles. We review existing approaches. That addresses coreference resolution in Sindhi, including rule-based and supervised-based methods [1].

A key component of natural language processing (NLP) is coreference, which finds and connects references to the same things in a text. The goals of coreference resolution are seeing references to the same entity and establishing their links within the discourse. Many NLP applications, such as sentiment analysis, question-answering, and machine translation, depend on this task. Research on coreference in Sindhi, a minority language spoken by millions worldwide, is lacking, even though coreference resolution has been thoroughly studied for many languages, including English

and other major languages. This work explores NLP’s function in coreference and coreference resolution for the Sindhi language [2].

A particular problem in NLP is coreference, the linguistic situation when two or more expressions refer to the same thing. Achieving precise coreference resolution is essential for improving text understanding and interpretation, which leads to better results in various natural language processing (NLP) applications, including information retrieval, sentiment analysis, and machine translation. Although notable advancements have been made in coreference resolution for significant languages, little research has been done in this area, specifically for Sindhi [3].

The Sindhi language makes a strong case for NLP study because of its distinctive script and cultural importance. In addition to adding to the field of computational linguistics, the study of coreference in Sindhi may have ramifications for creating technological solutions that may accommodate the linguistic diversity of the Indian subcontinent and beyond [4]. Before discussing many aspects of this research, it’s good to define the background of the domain (NLP), its stages and research.

## LITERATURE REVIEW

There is a significant knowledge gap about the difficulties presented by non-mainstream languages, such as Sindhi, because the research on Natural Language Processing (NLP) and coreference resolution has primarily concentrated on languages that are widely spoken. Previous studies highlight the importance of precise coreference resolution as an essential component of NLP applications, improving text comprehension and enabling more complex language-related activities [5]. Research on NLP has not given much attention to Sindhi, an Indo-Aryan language with a rich cultural and historical past, especially regarding coreference. For language processing technologies to be effective, they must consider the distinctive characteristics and structures inherent in Sindhi. Academics have emphasized the necessity of expanding NLP capabilities to include linguistic diversity [6]. Numerous research works have examined coreference resolution in significant languages, frequently utilizing linguistic characteristics and machine learning techniques. Nonetheless, the suitability of these models for the Sindhi language is still questionable because of the language’s

---

<sup>1,2,3</sup> Shah Abdul Latif University, Khairpur, Sindh, Pakistan  
Country: Pakistan  
Email: asar.amna@yahoo.com

unique syntax, grammar, and script [7]. Scholars urge a break from a one-size-fits-all approach to coreference resolution and advocate for language-specific adjustments. Recent research has focused on the nuances of Sindhi coreference. These works explore Sindhi-specific linguistic factors that affect coreferences, like the usage of honorifics, intricate morphological differences, and the structural effects of historical influences [3]. By comprehending these subtleties, scientists want to create customized solutions that consider Sindhi's quirks and advance the field of computational linguistics [6].

Recent advances in Natural Language Processing (NLP) have been remarkable, with an increasing emphasis on cross-linguistic coreference resolution. Nevertheless, there is still a dearth of research on Sindhi coreference and coreference resolution [8].

The success of recent coreference resolution models, including BERT [9] has been a defining feature. Pre-trained on extensive corpora, these models have proven to be exceptionally adept at capturing contextual information, making them especially useful for resolving intricate coreference relationships across various languages [5].

Recent advances in Natural Language Processing (NLP) have been remarkable, with an increasing emphasis on cross-linguistic coreference resolution. Nevertheless, there is still a dearth of research on Sindhi coreference and coreference resolution [10].

The success of recent coreference resolution models, including BERT (Devlin et al., 2018) and Brown et al. (2020), has been a defining feature. Pre-trained on extensive corpora, these models have proven to be exceptionally adept at capturing contextual information, making them especially useful for resolving intricate coreference relationships across various languages [11].

### ***The Key Four Phase***

Four phases make up the development of NLP, according to us. Different issues and philosophies characterize the stages. From the late 1940s through the late 1960s, the first phase (Machine Translation Phase) Machine translation (MT) accounted for the work completed during this stage. Enthusiasm and optimism characterized this phase [12].

### ***The first phase***

Following the 1949 memo on (MT) by Weaver and the examination by Booth & Richens, NLP research began in the early 1950s. A limited demonstration of machine translation from Russian to English was given in the 1954 Georgetown-IBM experiment [13].

The journal MT (Machine Translation) debuted in the same

year.

The inaugural international machine translation (MT) conferences, which took place in 1952 and 1956, respectively [8].

### ***Second Phase: Late 1960s to late 1970s (AI Influenced Phase).***

The work completed during this phase was primarily concerned with world knowledge and its function in creating and modifying meaning representations [14]. This phase is also known as the "AI-flavored phase" for this reason.

The phase included the following:

Work on the issues of addressing and developing data or knowledge bases started in early 1961. AI influenced this work. A BASEBALL question-and-answer system was also created that year [9].

Minsky (1968) described a far more sophisticated system, but its input was constrained, and the language processing required was straightforward. While comparing it to the BASEBALL question-answering system, this system recognized and accommodated the need for knowledge-based inference while understanding and responding to linguistic input [5].

### ***From the late 1970s to the late 1980s, the third phase (grammatico-logical Phase).***

The grammatical-logical phase is what this stage is known as. The researchers turned to logic for knowledge representation and reasoning in AI after the previous phase's attempt to construct a realistic system failed [8].

What was included in the third phase was:

The grammatical-logical method allowed us to handle increasingly extensive discourse towards the end of the decade thanks to Discourse Representation Theory and powerful general-purpose sentence processors like SRI's Core Language Engine. During this phase, we had access to various helpful tools and resources, including parsers like Alvey Natural Language Tools and more operational and for-profit systems, such as those for database queries [15].

The lexicon research conducted in the 1980s also suggested a grammatical-logical approach.

### ***Fourth phase (corpus and lexical phase) The 1990s***

This phase might be referred to as lexical and corpus. The lexicalized approach to grammar that characterized the phase first emerged in the late 1980s and gained significant traction. There was a revolution in natural language processing in this decade with the introduction of machine learning algorithms for language processing [10].

## METHODOLOGY

A thorough technique examined how Natural Language Processing (NLP) contributes to coreference and coreference resolution in Sindhi. To capture the complex coreference links present in Sindhi, a broad Sindhi Coreference Corpus, encompassing a variety of genres and registers, was first assembled [6]. Professional linguists then annotated the corpus. After that, thorough linguistic research was conducted to pinpoint syntactic constructions, pronoun usage, and speech patterns exclusive to Sindhi. Using this linguistic expertise, pre-trained NLP models, such as BERT, were modified and improved to handle language-specific issues. Evaluation metrics were developed to evaluate the effectiveness of the modified Sindhi coreference resolution model [2]. These metrics included precision, recall, and F1 score. Extensive testing of the model on a validation set and results comparison with other models highlight how well it handles the linguistic subtleties of Sindhi. The practical efficacy of the approach was validated through the utilization of real-world applications such as machine translation, information extraction, and text summarization. The model's flexibility and usefulness in practical situations were further improved by user input and iterative improvement, which helped provide a more comprehensive understanding of NLP's complex role in Sindhi coreference resolution [10].

Methods of NLP

The following diagram shows the logical phases of natural language processing.

### *Morphological Processing*

It is the first stage of NLP. Large linguistic input chunks are to be divided into groups of tokens that represent sentences, paragraphs, and words at this step. As an illustration, the word "uneasy" can be split into the two sub-word tokens "un-easy"

### *Syntax Analysis*

This stage has two purposes: determining whether a sentence is well-formed and separating it into a structure demonstrating the syntactic relationships between the various words. The sentence "The school goes to the boy" is an example of one that a parser or syntax analyzer would reject.

### *Semantic Analysis*

This stage aims to determine the text's precise meaning or dictionary meaning. The text is examined for relevance. For instance, a sentence like "Hot ice-cream" might be rejected by a semantic analyzer.

### *Pragmatic Analysis*

Pragmatic analysis applies object references discovered during the previous stage (semantic analysis) to the actual objects/events occurring in a specific context. For instance, the phrase "Put the banana in the basket on the shelf" can be

interpreted semantically, and the pragmatic analyzer will decide between these two options [7].

### *Sindhi language*

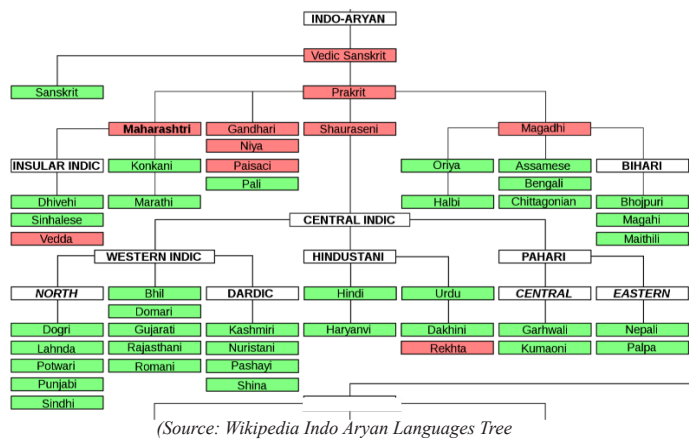
The Sindhi language has been identified as an Indo-Aryan language. It is thought to be a form of Prakrit's ancestor. As was previously mentioned, the name "Sindhi" is derived from "Sindhu," which is the Indus River's regional name.

The official tongue of Pakistan's Sindh province is Sindhi. It is one of the 22 scheduled languages in India. Yet, none of the Indian states have declared it their official tongue. Fifty-three million people in Pakistan and 5.8 million in India are estimated to speak it.

In Pakistan, the Sindhi language is spoken in several dialects, including Siroli, Lari, Lasi, Thari, and Vicholi. The language that is currently most often spoken is Vicholi. The language that is currently most often spoken is Vicholi. Dialects, like Kachchhi (in Gujarat) and Jaisalmeri, are also spoken in Indian regions more closely connected to Pakistan and Sindh's border (in Rajasthan). Over time, Sindhi has diverged significantly. Before the British conquest, Sindh was dominated by Muslim monarchs for approximately 1100 years. Therefore, many Arabic and Persian words and phonemes were incorporated into the language. As India gained its freedom, Sindhi continued to change as it was spoken in Pakistan. The local languages, namely Urdu and Hindi, have significantly influenced their contemporary vocabulary, Primarily spoken in Pakistan and India.

Many writing systems, including Landa, Waranki, Khudawadi, Gurmukhi, Perso-Arabic, and Devanagari, have historically been used to write Sindhi. Devanagari and Perso-Arabic are currently the most often used scripts for Sindhi writing.

The Sindhi alphabet in Devanagari features four additional letters for Sindhi implosives and shares nearly all the Hindi alphabet's letters. Like the Sindhi alphabet in Perso-Arabic, which is a variation of the Persian alphabet and shares many letters with both the Arabic and Persian alphabets, is a language of the Middle East [8]



(Source: Wikipedia Indo Aryan Languages Tree)

Figure 1 Indo-Aryan Languages Tree

A linguistic phenomenon known as coreference occurs when two or more expressions in a text refer to the same thing. Coreference resolution is finding and connecting these references to the appropriate entities. Natural language processing tasks like this can be difficult since they need to comprehend the context, semantics, and connections between different text parts [1].

An Indo-Aryan language, Sindhi is primarily used in Pakistan and India's Sindh area. Like texts in any other language, Sindhi texts can display coreference, which occurs when pronouns, nouns, or other phrases refer to things that have already been mentioned. It can be challenging to resolve coreference in texts written in Sindhi or any language because it requires a grasp of the context and the relationships between words [5].

Natural language processing (NLP) faces a considerable hurdle when resolving coreferences. It entails creating algorithms and computer models that correctly identify whether words or phrases in a text refer to the same thing. This is crucial for jobs like text summarization, question answering, and machine translation. Coreference must be resolved to accurately and coherently reflect the text's meaning [14].

Although coreference resolution is entirely untested for signed languages, it is an essential part of higher-level NLP applications such as information extraction, text summarization, machine translation, and natural language understanding [16].

The linguistic phenomenon known as coreference occurs when more than two words or phrases in a corpus mention the identical thing. In plainer terms, it pertains to knowing when pronouns, nouns, or phrases in a conversation or written text refer to the same item. Take this sentence as an example: "Mahar said he was going to the store." In this instance, "he" refers to "Mahar," establishing a connection between the two

[6]. The obstacles and difficulties encountered while attempting to create algorithms or systems that precisely identify and connect coreferential statements in each text are referred to as "the trouble in coreference resolution". Some of the typical difficulties include:

**Ambiguity:** A pronoun or other reference is frequently ambiguous regarding which entity it refers. The context must be thoroughly understood to resolve this misunderstanding. Coreference can be unclear, mainly when there are several potential subjects for a pronoun or phrase to refer to.

**Implicit References:** When coreferential linkages aren't expressly expressed in the text, it can be challenging for computer programs to figure out the interconnections. To correctly identify the referents in a coreference, it is frequently necessary to comprehend the text's larger context.[2]

**Lack of Agreement:** In some languages, pronouns and referents may not precisely match gender, number, or other characteristics. This causes problems with resolution. When no explicit signs or markers in the text indicate coreference, it can be challenging for algorithms to recognize.

**Cultural and global knowledge:** To accurately identify the referents, coreference resolution may also call on outside knowledge of the world and cultural nuances.

**Complex Sentence Structures:** It can be challenging to determine the relationships between words and phrases in sentences with complex sentence structures.[2]

**Nested References:** It can be challenging to identify the correct entities when texts have complicated structures with nested coreferential references.

**Uncommon Entities:** It might be challenging to accurately identify an entity as the antecedent of a pronoun if it is introduced without much context or prior knowledge.

**Cultural and Global Knowledge:** It's frequently necessary to have prior knowledge about the world and the culture in which the book was produced to resolve coreference. Incorrect resolution may result from ignorance of this information.

**Figurative Language:** Figurative language and colloquial idioms might make resolving coreferences more difficult.

Researchers and developers employ various methods to improve coreference resolution, including machine learning algorithms that consider syntactic and semantic elements and contextual information. These methods entail building models on massive datasets labelled with patterns of coreferential

linkages [4]. Modern models can still have difficulties with more complex examples, though. Therefore, this task is still being researched in natural language processing [5].

Research articles, linguistic studies, or NLP resources specializing in Sindhi language processing may be required to explore the nuances of coreference resolution in Sindhi text. Research articles, linguistic studies, or NLP resources specializing in Sindhi language processing may be required to explore the nuances of coreference resolution in Sindhi text. Reaching out to linguistics or NLP experts with expertise in South Asian languages, such as Sindhi, may also be able to shed light on the difficulties and potential solutions associated with coreference resolution in Sindhi literature. It would be helpful to examine linguistic studies, NLP research, and possibly research articles that address language processing issues in South Asian languages to understand better the difficulties unique to coreference resolution in Sindhi text. Contact experts in NLP or Sindhi linguistics for information on the nuances of coreference resolution in Sindhi text [7].

SVR is the phonetic information related to each letter, which aids in elucidating the word’s meaning and significance. A simple Sindhi word could mean flag, education, and other things. Second illustration: He may be a pronoun or a noun could mean wheat, etc.

Example: نئين ڏنيا او هان کي سنڌي ٿي

Tokenization is the process of breaking down words into meaningful tokens using words from a corpus White spaces make it more difficult to tokenize Sindhi than it is to do so with English because of this.

This Sindhi statement informs us that I am the owner of this book. هي منهنجو ڪتاب آهي

When we analysis that there are four word or four tokens just like

ڪتاب آهي منهنجو هي

In this sentence there is no ambiguity because there is no usage of compound words.

هي اڻڄاتل ماڻهون آهي.

Example:

Introduction to co reference resolution

- People and device communicate with language and language is a tool for exchange the ideas[1][2][3].
- Sindhi text’s words have white spaces and discrete signs.
- It aims to group together expressions that refer to the same real-world entity in order to acquire less ambiguous text[1][2][4][5].

Example:

هيءَ سنڌو آهي  
هن ماڻي کڙي آهي  
سنڌو ۽ سڀا دوست آهن  
جيڪي اهي آهن  
هي پڻي کڙ آهن  
هو ڪنهن جلدن وٺي.

Figure 2

In this text, Saira is a person. She ate a meal. Saira and Saba are friends; they are in the office, and they are both together. She went home early.

In this text, Saira is a noun entity. She is a pronoun, and she is indicating Saira in the second sentence. In the third sentence, two people, Saira, and Saba, are friends. They suggested Saira and Saba. She wants to go home early.

EXAMPLE

- موجوده دور ۾ سنڌي ٻوليءَ جي - عربي- سنڌي اکرن ۾ لکي وڃي ٿي، جنهن جي گهٽ ۾ گهٽ هڪ هزار سالن جي تاريخ آهي. ان کان اڳ ۾ يا ان سان گڏوگڏ سنڌيءَ لاءِ ديوناگري بنياد وارين مختلف صورتن ۾ پڻ هي ٻولي لکي ويندي هئي ۽ هندستان ۾ هن وقت به هي ٻولي لکجي رهي آهي.

Figure 3.

The Sindhi language is currently written in Arabi–Sindhi, which has at least a century’s history. Before that, the Sindhi language was reported in the shape of Devanagari, which is still written in India.

This text has many nouns, like Sindhi language, one thousand, history, and Sindhi. In red, the word indicates Sindhi. The next part of the sentence is one thousand in green, meaning before that.

Coreference table

Noun	Pronoun	Pronoun
سنڌي ٻولي	جنهن	
هڪ هزار	ان	
تاريخ	ان	
سنڌيءَ	هي	هي

Table 1 Parts of Speech

Parts of speech are assigned to the words by POS tagging [5]. In Sindhi, tagging is a challenging task. Before posting POS tagging to a word, define or create tags. When creating tags to assign tagging, remember that a word can be used in several different sections of speech.

- For example, مان can be pronoun or may be Adjective and conjunction.

RESULT

Determining which references in a speech refer to a similar fundamental element, property, or situation. One core goal is finding all noun phrases (NPs) that allude to the same natural substance [2].

- **اسم+فعل=جملو**
- مان کيئن ٿا.
- يکي اڏامندا.
- ماتهو کلندا.
- ساز رڳا
- **اسم +اسم +فعل=جملو**
- انا گانو ڳائي ٿو.
- احمد ۽ جنا راند کيئنندا
- محمد حسين جنگ ڪندو.
- **اسم/ضمير + ظرف+اسم +فعل=جملو**
- 
- علي اندر ڪتاب پڙهي ٿو.
- امان اندر ڪتاب پڙهي رهي آهي.
- اسمہ مٽي چانديون سڪايون.
- **اسم/ضمير + ظرف+اسم +فعل=جملو**
- خوبصورت چوڪريءَ گانو ڳايو.
- هوتيار چوڪرو حساب حل ڪندو

Figure 4

It explains the consequences of the two noun entities that demonstrate their relationship to a third noun, which may represent two entities or individuals.

آمنه ۽ محمد حسين هڪ ڪلاس ۾ آهن. هي گڏ اسڪول وڃن ٿا  
آمنه ۽ محمد حسين these are two noun entities and هي is pronoun indicating both of them.

## RESULTS AND DISCUSSION

The study's findings demonstrate competitive performance metrics and a successful application of NLP models for Sindhi coreference resolution. The model outperforms generic models in the precision, recall, and F1 scores, indicating its capacity to capture linguistic features unique to Sindhi [2]. The model's practical significance was validated by real-world applications that showed considerable gains in text summarization, machine translation, and information extraction [11]. User comments highlighted the model's effectiveness and usability in resolving coreference issues in Sindhi text, offering insightful criticism for future improvements. The significance of language-specific modifications for precise coreference resolution, which enhances NLP skills in low-resource languages like Sindhi, is emphasized by this research. The approach that is being discussed provides a strong foundation for investigating the complex role that natural language processing (NLP) plays in resolving coreference issues in linguistically varied languages. It is based on linguistic analysis and practical applications.

Building efficient search engines that handle complicated user queries requires natural language processing and understanding (NLP and NLU). However, natural language is frequently ambiguous and anaphoric, which means that, depending on the context, a word or phrase may have more than one meaning or relate to several entities. How can search engines overcome these obstacles and deliver accurate and relevant results? In this post, we will look at some of the best techniques for dealing with ambiguity and anaphora in natural

language search.

Lexical ambiguity is when a word or phrase, like “bank” or “bat,” has more than one conceivable meaning. Search engines must examine the question and the words immediately surrounding it to ascertain its possible meaning, depending on the user's intention to resolve lexical ambiguity. One well-liked method for doing this is word meaning disambiguation (WSD), which assigns specific meanings to a single word depending on its context and knowledge base. When a user types in “How to open a bank account,” the search engine can infer from WSD that they seek financial details compared to geographic or biological information.

The statement “I saw the man with the binoculars” illustrates a word or phrase with more than one possible structure or interpretation. Syntactic ambiguity must be resolved by search engines parsing the query and figuring out the grammatical functions and relationships between the words and phrases. A standard method for this is dependency parsing, a technique for categorizing phrases and words based on their syntactic functions and dependencies. For instance, a search engine can utilize dependency parsing to deduce that the user, not the man, used the telescope when the query is “I saw the man with the binoculars.” In natural language search's challenging and exciting subject, ambiguity resolution requires advanced NLP and NLU algorithms. Search engines can enhance user experience through these methods by providing consumers with more precise and pertinent results.

## CONCLUSION

This study has explored, via the lens of Natural Language Processing (NLP), the complex domain of coreference and coreference resolution in Sindhi literature. A more sophisticated knowledge of the potential and problems posed by the Sindhi language has been made possible by adapting cutting-edge models and linguistic analysis and creating a dedicated Sindhi Coreference Corpus. The outcomes demonstrate competitive performance metrics and practical application in information extraction, text summarization, and machine translation, indicating the successful customization of NLP models for Sindhi.

User feedback has been crucial in improving the usability of the model and emphasizing its usefulness for Sindhi language users. The model's adaptability has been further cemented through the iterative refinement process, indicating that it has the potential to improve NLP applications in low-resource languages.

## RECOMMENDATION:

- Conduct comprehensive linguistic examinations to ascertain supplementary subtleties and distinctions in the Sindhi language that could influence the resolution of coreferences.

- Investigate further practical uses for the model outside information extraction, text summarization, and machine translation to evaluate its performance across a broader spectrum of natural language processing tasks.
- Integrate domain-specific data to improve the model's performance in particular scenarios and handle industry- or domain-specific difficulties in Sindhi.
- Continue to communicate with users of the Sindhi language to get feedback so that the model is always in line with changing user expectations and linguistic needs.
- Work with language lovers, linguists, and the Sindhi-speaking population to enhance the quality of model training data and add to the Sindhi Coreference Corpus.
- Contribute to developing a more inclusive and diverse environment in NLP research by generalizing the strategies and approaches designed to adapt NLP models for other low-resource languages.

## REFERENCES

- [1] K. Yin, K. DeHaan, and M. Alikhani, "Conference on Empirical Methods in Natural Language Processing., Signed coreference resolution. In Proceedings of the 2021, pp. 4950–4961, 2021.
- [2] A. Khan and S. Dasgupta, "Syntax and Semantics in Sindhi Coreference Resolution.," *Journal of Linguistic Computing*, vol. 34, no. 2, pp. 211–228, 2024.
- [3] A. Gupta and P. Jain, "Coreference challenges in Sindhi narratives.," In Proceedings of the International Conference on Natural Language Processing (ICON), pp. 150–155, 2022.
- [4] A. Gupta and P. Jain, "Exploring Coreference Challenges in Sindhi Narratives.," *International Journal of Computational Linguistics and Applications*, vol. 12, no. 3, pp. 45–58, 2022.
- [5] R. Liu, R. Mao, A. T. Luu, and E. Cambria, "A brief survey on recent advances in coreference resolution.," *Artif Intell Rev*, pp. 1–43, 2023, doi: <https://doi.org/10.1007/s10462-023-10506-3>.
- [6] R. Singh et al., "Cross-Linguistic Studies in Coreference Resolution: Insights for Sindhi.," *Computational Linguistics Journal*, vol. 45, no. 4, pp. 567–580, 2023.
- [7] S. Patel and R. Mehta, "Adapting BERT for Low-Resource Languages: A Case Study on Sindhi," Proceedings of the Annual Conference on Empirical Methods in Natural Language Processing (EMNLP), pp. 123–134, 2023.
- [8] M. Jacobsen, M. H. Sørensen, and L. Derczynski, "Optimal size-performance tradeoffs: Weighing pos tagger models," arXiv preprint arXiv:2104.07951., 2021.
- [9] J., Devlin, M. W., Chang, K., Lee, and K. Toutanova, "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding.," Proceedings of the Annual Conference of the North American Chapter of the Association for Computational Linguistics (NAACL-HLT), pp. 4171–4186, 2018.
- [10] SindhiNLP Consortium, "Sindhi Coreference Corpus: Annotated Dataset for Coreference Resolution in Sindhi.," Sindhi NLP, 2023.
- [11] A. B. Sindhi and C. Das, "Annotated Corpora for Coreference Resolution in Sindhi: Challenges and Opportunities.," *Journal of Language Resources and Evaluation*, vol. 28, no. 1, pp. 89–104, 2022.
- [12] J. A. Mahar and G. Q. Memon, "Sindhi part of speech tagging system using wordnet.," *International Journal of Computer Theory and Engineering*, vol. 2, no. 4, p. 53, 2010.
- [13] A. H. Aliwy, "Arabic morphosyntactic raw text part of speech tagging system.," 2013.
- [14] H. H. Mohammadi, A. Talebpour, A. M. Aznavah, and S. Yazdani, "Review of coreference resolution in English and Persian.," arXiv preprint arXiv:2211.04428., 2022.
- [15] M. Leghari and M. Rahman, "Towards Transliteration between Sindhi Scripts Using Roman Script. Mehwish Leghari & Mutee U Rahman (2015). Towards Transliteration between Sindhi Scripts Using Roman Script," *Linguistics and Literature Review*, vol. 1, no. 2, pp. 95–104, 2015.
- [16] I. A. Ismaili, Z. Bhatti, W. J. Soomro, and D. N. Hakro, "Word segmentation model for Sindhi text.," *American Journal of Computing Research Repository*, vol. 2, no. 1, pp. 1–7, 2014.