# RULE-BASED IMPLEMENTATION OF GUJARATI SANDHI DISAMBIGUATION USING NLP

**Mr. Nitesh G. Patel**, Ph.D. Scholar, Department Of Computer Science, Gujarat Vidyapith, Ahmadabad, Gujarat :: nitesh.mscit@gmail.com
**Dr. Dhiren B. Patel**, Professor, Department Of Computer Science, Gujarat Vidyapith, Ahmadabad, Gujarat dhiren_b_patel@gujaratvidyapith.org

**Abstract:**
Sandhi refers to the merging of letters or sounds at word boundaries, causing difficulties in text processing. In other words, Sandhi is a phenomenon where letters or sounds at word boundaries mix, changing pronunciation and shape, which challenges Gujarati text processing. This study recommends a rule-based approach to handling sandhi transformations. This paper tackles the challenge of sandhi disambiguation in Gujarati Natural Language Processing (NLP). This paper proposes a rule-based framework that uses AI and Machine Learning techniques for sandhi disambiguation in Gujarati NLP. We describe the architecture, the system's implementation, and assessment measures, including processing time and accuracy. We inspect how this method enhances different NLP tasks. We report on task performance with sandhi handling part-of-speech tagging accuracy, and text segmentation improvement. We also look into how the sandhi-handling approach may be used in machine translation. We assess its effect on the implementation of Gujarati sandhi's accuracy using several indicators & outcomes of human review. It is expected that the outcomes will show how well the rule-based method handles sandhi, producing remarkable gains in NLP tasks in Gujarati grammar.

**Keywords:** Gujarati, Machine learning, rule-based NLP, sandhi, segmentation

**Introduction:**
This Research explores the potential of Natural Language Processing (NLP) techniques for implementing rule-based Gujarati grammar. Gujarati, a prominent Indian language, possesses unique grammatical features like sandhi (euphonic combinations) that pretense challenges for NLP applications. This work aims to bridge this gap by outlining a framework for incorporating Gujarati grammar rules into NLP systems.

It's crucial to understand that Gujarati is more than just Sanskrit. Gujarati has also incorporated elements of Arabic, Persian, and English throughout the millennia. Gujarati is a distinctive and lively language because of its diverse range of influences. There are historical and cultural ties between Gujarati and Sanskrit. Gaining insight into this relationship helps one to appreciate the Gujarati language's development and history on a deeper level. It is akin to watch upon a family tree from which Gujarati proudly descends as a branch from the great Sanskrit hierarchy.

Basic-Well known concepts of Gujarati Grammar are: કૃદંત,નિપાત,વિભક્તિ, સંધિ,સમાસ,સમાનાર્થી શબ્દ,જોડણી,અલંકાર,છંદ,લિંગ અને વચન,વિરુદ્ધાર્થી શબ્દો, શબ્દ સમૂહ માટે એક શબ્દ,રૂઢિપ્રયોગ,કહેવત.

Every concept has its own rules and distinct logic. Like in other languages, the majority of the rules are extremely rigorous and unambiguous. These ideas can be put into practice on a computer by using certain morphological rules. This study's main goal is to assess the state of the art in terms of technology and the best way to use specific NLP techniques to implement the Sandhi notion of Gujarati grammar. Because Gujarati, like many Indian scripting languages, has its own set of rules for mixing consonants, vowels, and modifiers, our main focus is on the rule-based implementation.
In Gujarati, Sandhi:

The term "sandhi" in Gujarati refers to the phonetic combination of words in the language with clearly defined rules. The word "sandhi" (Sama + Dhi) means "similarity" or "joint." In Sanskrit, Hindi, and other languages, the disorder produced due to the combination of interchangeable vowels or letters is called "sandhi." Sandhi leads to phonetic transformation at word boundaries of a written chunk (small part), and the sounds at the end of a word join together to form a single chunk of the character sequence.

- Some Gujarati Sandhi Examples :

નિસ્ + ચિંત = નિશ્ચિત,          અંતઃ  + કરણ = અંતઃકરણ

સમ્+ તોષ = સંતોષ,          સપ્ત  + ઋષિ = સપ્તર્ષિ

In Gujarati Language there are majorly 4 Types of sandhi is Present:

1. સ્વર સંધિ : અ,આ,ઇ,ઈ,ઉ,ઊ,ઋ,એ,ઐ,ઓ,ઔ

2. વ્યંજનસંધિ:

ક,ખ,ગ,ઘ,ચ,છ,જ,ઝ,ઞ,ટ,ઠ,ડ,ઢ,ણ,ત,થ,દ,ધ,ન,પ,ફ,બ,ભ,મ,ય,ર,લ,ળ,વ,શ,ષ,સ,હ,ૂ

3. વિસર્ગ સંધિ :   ":" based

4. અનુનાસિક સંધિ

We are majorly focused on Swar sandhi (સ્વર સંધિ) in starting phase as it is the initial and fundamental clear-cut rules based Sandhi.

(સ્વર સંધિ) Swara Sandhi also distributed in other 5 parts (As per Sanskrit rules):

1. Dirgha Sandhi (લાંબા સંધિ),  2. Gun Sandhi (ગુણ સંધિ), 3. Vriddhi Sandhi (વૃદ્ધિ સંધિ),  4. Yan Sandhi ( યાન સંધિ ), 5. Ayadi Sandhi (આયદીસંધિ)


**Concept of Gujarati words before proceed to Sandhi:**

Before proceed towards Sandhi concept we have to Take care of some of the grammatical concepts: As, we can filter Words before sandhi Vichada using concept given below:

The most important and special thing to remember about the concept of Sandhi vigraha is not originate from Gujarati language. The Gujarati language has a large vocabulary, and from this large vocabulary only one class of words related words adopted from Sanskrit into Gujarati has traditionally undergone word convention.

The same rules apply in Gujarati language as in Sanskrit language for syndication of these (તત્સમ) Tatsama words. Tatsama words from non-Sanskrit languages or (તદ્ભવ) Tadbhava words from other languages including Sanskrit or local indigenous words never have Sandhi-Vigraha.

The Concept come in to the reorganization that is called: (તત્સમ & તદ્ભવ શબ્દો)

1. તત્સમ શબ્દો: આ શબ્દો સીધા સંસ્કૃત ભાષામાંથી આવે છે.

e.g.: શિક્ષક (shikshak) - teacher, ધર્મ (dharma) - religion

Other e.g.: સૂર્ય, ચંદ્ર, મયુર, અક્ષર , વિદ્યુત , શૂન્ય, હસ્ત

2. તદ્ભવ શબ્દો: આ શબ્દો સંસ્કૃત ભાષામાંથી પરિવર્તન થઈ ને આવે છે.

e.g.: મનનું (mannuṁ) – mind, સાચું (sāchuṁ) - true

Other e.g.: ખીર, સખી, આંખ , હાથ , સસરા, વારસ

From them we can Perform Sandhi Only on "તત્સમ" words as they are directly taken from Sanskrit Language.

**Review of Literature:**

In 2011, Nair and Peter developed a rule-based system for splitting compound words in Malayalam, enhancing language processing techniques.[9] In 2016, Joshi and Kushwah presented an algorithm for Sandhi in Hindi, achieving approximately 98% accuracy in compound word formation .[5] Cao and Rei (2016) introduced a joint model for word embeddings and morphology, addressing limitations in traditional models .[3] Krishna et al. (2016) built a word segmenter for Sanskrit, showcasing advancements in language processing methodologies .Researchers reviewed South Asian languages, highlighting linguistic diversity and challenges in language processing.[2]

Patel and Patel (2018) reviewed rule-based Gujarati grammar implementation, emphasizing morphological rules and text processing advancements.[11] In 2020, Yadav and Kumar explored NLP applications in Indian languages, focusing on morphological analysis and its significance .[10] In 2021,

Modh and Saini focused on idiom detection in Gujarati, utilizing diacritics and suffix-based rules for improved accuracy , they also collected 3,240 unique n-gram Gujarati idioms, enhancing idiom detection capabilities .[21]Baxi and Bhatt (2024) developed a bidirectional LSTM-based morphological analyzer for Gujarati, achieving significant accuracy improvements in processing.[1]

Dave and Mehta (2023) evaluated stemmers for Gujarati, finding hybrid approaches most effective across various datasets.[15] Kevat and Degadwala (2023) systematically reviewed Gujarati-text summarization methodologies, emphasizing the need for diverse datasets.[6] Rathi and Gupta (2022) studied hybrid approaches in Gujarati text summarization, contributing to the field's understanding of summarization techniques.[13]

Thakkar and Mehta (2022) provided a comprehensive survey on text summarization techniques for Gujarati language processing.[20] In 2023, Zaveri and Shah reviewed sentiment analysis techniques for Gujarati, identifying gaps in existing research and methodologies.[22] The literature emphasizes the importance of linguistic components in enhancing text processing for Indian languages .[19]

Research highlights the effectiveness of extractive summarization techniques, particularly for Gujarati text .The studies collectively advance the understanding of Gujarati and other languages processing, addressing unique linguistic challenges. Overall, these contributions pave the way for effective grammar implementation in Indian NLP applications.

**Problem statement:**

The most important and interesting challenge of sandhi implementation is to generate a system or facility for common people who are using the of Gujarati language. Millions of people speak Gujarati in the whole world, however there are not enough refined computing resources or tools to analyze and understand its grammar using natural language processing (NLP) approaches. This gap offers a special opportunity for research and development as well as a distinct challenge.

The main motivation of this research is to create a reliable, rule-based system that uses natural language processing (NLP) to implement Gujarati grammar. By utilizing rule-based approaches, the system seeks to address numerous important facets of Gujarati grammar, such as syntax, morphology, and semantics. Using this method, a thorough set of grammatical rules will be developed that can be utilized to both evaluate and produce grammatically accurate Gujarati writing.

Important issues that this research needs to resolve are Gujarati Grammar's Complexity, Insufficient Resources for Utilizing Current NLP Tools, Performance Efficiency and Precision. Our primary goal is to use ML and NLP principles to execute "sandhi" in a rule-based manner.
There are two primary methods exist for putting "sandhi" into execution:

(1) Joining of Sandhi (સંધિ છોડો) & (2) Sandhi Vicched (સંધિ જોડો)

There are main 4 Types Of sandhi and they have their individual rules set. Among that we have focused on the Swar Sandhi.

| સ્વર સંધિ : | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| **Swar Sandhi Rules** | | | **Example Set :1** | | | **Example Set :2** | | |
| અ | અ | આ | સૂર્ય | અસ્ત | સૂર્યાસ્ત | પરમ | અર્થ | પરમાર્થ |

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| ઇ | ઇ/ઈ | ઈ | પરિ | ઈક્ષા | પરીક્ષા | રવિ | ઇન્દ્ર | રવીન્દ્ર |
| ઉ | ઉ/ઊ | ઊ | ભાનુ | ઉદય | ભાનૂદય | સિંધુ | ઊર્મિ | સિંધુર્મી |
| અ/આ | ઇ | એ | ઉપ | ઇન્દ્ર | ઉપેન્દ્ર | યથા | ઈસ્ટ | યથેષ્ટ |
| અ/આ | ઉ/ઊ | ઓ | સૂર્ય | ઉદય | સૂર્યોદય | ગંગા | ઊર્મિ | ગંગોર્મિ |
| અ | ઋ | અર્ | સપ્ત | ઋષિ | સપ્તર્ષિ | મહા | ઋષિ | મહર્ષિ |
| અ/આ | એ | ઐ | એક | એક | એકૈક | સદા | એવ | સદૈવ |
| અ/આ | ઔ | ઔ | વન | ઔષધિ | વનૌષધિ | મહા | ઔષધિ | મહૌષધિ |
| ઇ | વિજાતીય | ય્ | પ્રતિ | એક | પ્રત્યેક | અતિ | આનંદ | અત્યાનંદ |
| ઉ/ઊ | વિજાતીય | વ્ | સુ | આગત | સ્વાગત | પૃથુ | ઇ | પૃથ્વી |
| ઋ | વિજાતીય | ર્ | પિતૃ | આજ્ઞા | પિત્રાજ્ઞા | માતૃ | આદેશ | માત્રાદેશ |
| એ | વિજાતીય | અય | ને | અન | નયન | મે | ઊર | મયુર |
| ઐ | વિજાતીય | આય | નૈ | ઈકા | નાયિકા | ગે | અન | ગાયન |
| ઔ | વિજાતીય | આવ | નૌ | ઈક | નાવિક | પૌ | અન | પાવન |
| ઓ | વિજાતીય | અવ | પો | અન | પવન | ભો | અન | ભવન |

**[Table : Swar Sandhi Rules with Suitable Example ]**

The Swar sandhi based on the total 11-swar (અ, આ, ઇ, ઈ, ઉ, ઊ, એ, ઓ, ઐ, ઔ, ઋ) of the Gujarati phonetics. From the below rule set that is also cleared that there are two part of these rules also that is: સજાતીય સ્વર & વિજાતીય સ્વર.

**Research Objectives:**

- Using NLP approaches construct an absolute rule-based Gujarati grammar system.
- The practical implementation of "sandhi (સંધિ)" concept of Gujarati Grammar in mutual ways: Join multiple words to generate a single sandhi word & Separate the multiple words from single word sandhi
- Collect the Gujarati Language's words from authorized source (like: Gujarati Lexicon).
- Accumulate the proper Gujarati grammatical rules (Morphological Rules) that addresses syntax, morphology, and semantics.
- Verify that the system works with the NLP frameworks and tools that are currently in use.
- Manage both the databases one is for Gujarati Words for spell checking & another is sandhi words for testing purpose.
- Conduct thorough testing and validation to assess the developed system's accuracy and performance.

**Data Collection:** The data (Gujarati language words) is received from Gujarati Lexicon.

Details of data: Total words: 12949 :(Words in consolidate excel file) : 'ક', 'ક્ષ', ' ખ', 'ગ'.

- Details of data (words in excel file) :
  - 7000 words starting from letter  : 'ક'
  - 180 words starting from  letter   : 'ક્ષ'
  - 2495 words starting from letter  : ' ખ'
  - 3170 words starting from letter  : 'ગ'

We can take words input for data from any type of text box or Tool which is available in from of Gujarati Language or we can take Gujarati Sentences & after that we can Split the Words by Space. Store the words in Unicode Gujarati wise in Database. Check for the correct Spelling & Classify the word in any of the category which is already discussed.

Another file generated by me for well-known available sandhi (Gujarati Grammar) combinations. It consist of 210+ sandhi combinations (languidly correct sandhi combination) with all 3 types available types of sandhi. These data is collected from various websites, books, online videos & other known resources. This data is very useful for checking or testing of data while implementation of sandhi.

After receiving the words the first step is to clean the words which are received. Cleaning the words Process involves steps like:

– Remove the stop words, Remove the complex Multiple words, Remove the Multiple Repetition Words ,Remove the single Character words

**Research methodology:**

Words Classification: Distributed Words Based on different Classes & categories displayed below.

**Class Category 1: (Eligible for Sandhi)**

- **CLASS   A :** Tatsam Words Directly Borrowed from Sanskrit
- **CLASS   B :** Swar based Sandhi Words
- **CLASS   C :** Vyanjan based Sandhi Words
- **CLASS   D :** Anunashik based Sandhi Words
- **CLASS   E :** Visharga based Sandhi Words

**Class Category 2: (Not-Eligible for Sandhi)**

- **CLASS   F :** Stop words (Not eligible for Sandhi)
- **CLASS   G :** Tadbhav words (Not eligible for Sandhi)
- **CLASS   H :** Single Character Word (Not eligible for Sandhi)
- **CLASS   I   :** Other Special Words from other Languages
- **CLASS   J :** Words Originated from Gujarati itself.

**Sandhi Classification:**  This above classification is based on Gujarati grammar's Relevant Rules from Books & Other References.

1. Swar Sandhi ( વિદ્યા + આલય = વિદ્યાલય )

2. Vyanjan Sandhi ( સમ્ + હાર = સંહાર )

3. Visarga Sandhi ( મનઃ + રથ = મનોરથ )
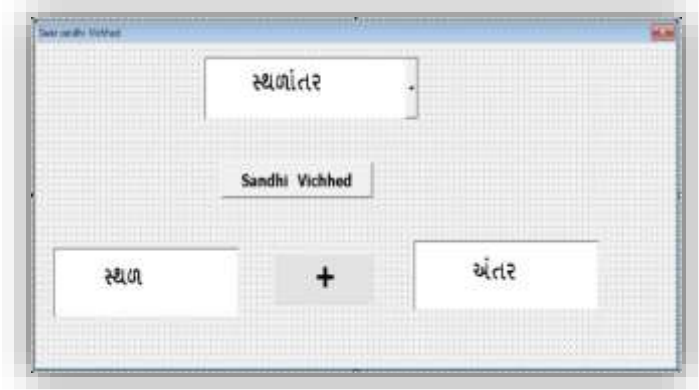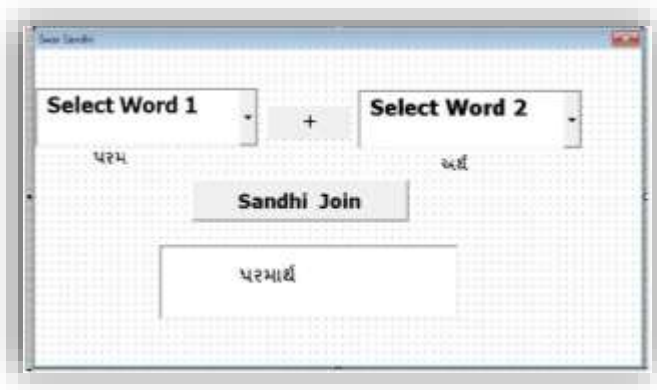
4. Anu-Nashik Sandhi ( સત્ + નારી =સન્નારી )

**Basic Level Sandhi Implementation:**

Before Actual implementation of Sandhi with AI &ML, we have developed one small tool in VB. NET. All the data feed to the tool is stored in Excel  Sheet According to rules the sandhi joining & Sandhi Splitting is performed in Visual Basic script.

It is simply perform sandhi joining as 1ˢᵗ words' list with 2ⁿᵈ words' list; if it is suitable as per rule then sandhi joining will perform. Same vice versa we can choose one word from list and if it eligible for sandhi Vicched then tool will display the sandhi words. This Tool Consist of two Parts:
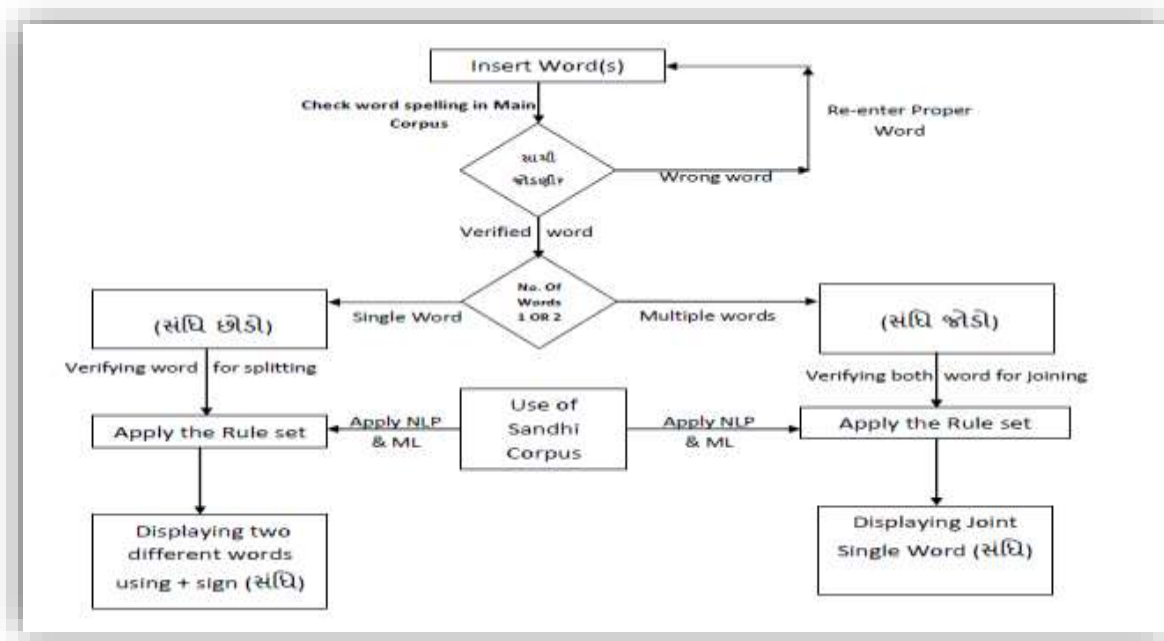
**1. Sandhi Joining**                          **2. Sandhi Splitting (Vicched)**



[Fig.: Sandhi Joining & Splitting Tools]

**Sandhi Implementation Approach flow-chart:**



[Fig.: Sandhi Implementation Approach flow-chart]

From the available two methodologies we concentration on:  Separate the multiple words from single word sandhi. As per display in the diagram we can approach to the actual implementation of Gujarati grammar's "sandhi" stepwise as demonstrate below:

**Step-1:** Initially, we have to enter the word manually into the system or we can copy-paste the word into the input box or facility.

**Step-2:** We have to check whether the inputted or copied word is available in our language corpus.

**Step-3:** If the same (whole) word is not available in the corpus then we have to take input again. If the word is available in corpus then we can go ahead.

**Step-4:** Now we have to check the word is split-able (sandhi) or not. So, we have to apply some of the concepts of machine learning and some NLP rules for this type of checking.

**Step-5:** Each word of Gujarati language is not split-able using sandhi rules and if we impose rules forcibly then the grammar it-self will go wrong.
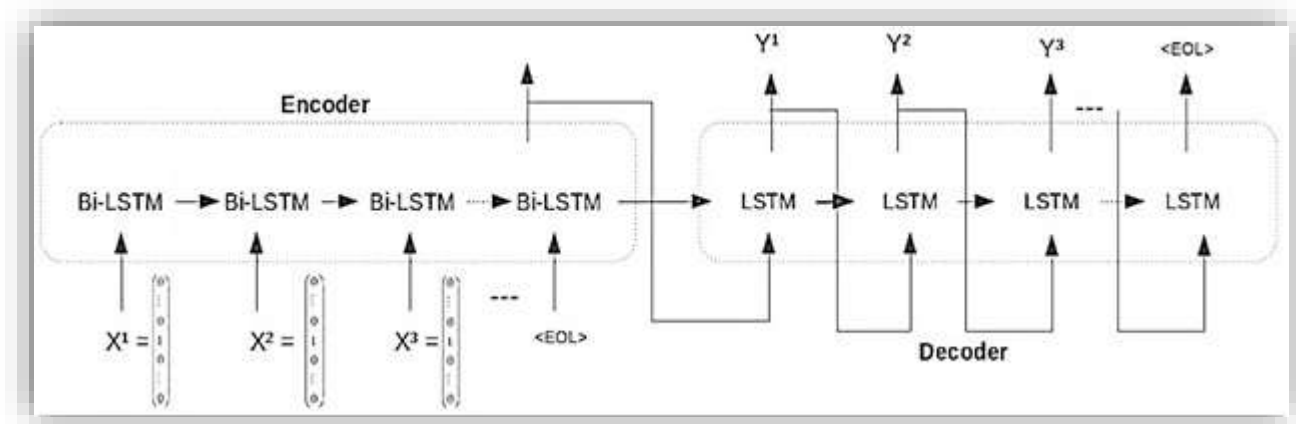
**Step-6:** So, it's clear that each word of the Gujarati dictionary is not compatible with "sandhi". After this most important step, we have to consider which type of sandhi rule we can apply (among all 4 Types) on the word.

**Step-7:** We have to apply some rule-based NLP concepts to implement at this stage.

**Step-8:** Finally, after execution, we can get the proper sandhi words from our single inputted word.

If we follow the above process in a reversible manner then maybe we can get proper Gujarati word, which is available in the corpus that is our final goal.

**Model Architecture used for Sandhi:**



[Fig.: NN Architecture for Generating Sandhi]

The model is trained on a dataset of Gujarati compound words and their corresponding split words. During training, the model learns to map the input word to its split word representation. Once the model is trained, it can be used to generate split words for new input words. Overall, the block diagram you sent me shows a neural network architecture that can be used to generate and split compound words in Gujarati.

This architecture is based on bi-directional LSTM networks and attention mechanism, which are both powerful techniques for natural language processing tasks.

**Stage 1** - Predict the sandhi-window in a compound word that is to be split, using a RNN model. All elements in this array are 0 except the elements corresponding to the sandhi-window characters which are set to 1.

**Stage 2** - The sandhi-window is then split in 2 words using a seq2seq model. Let's call the first split of sandhi-window as $S1$ and second split of sandhi-window as $S2$.

**Post-Processing step** - The characters $NW1$ before the sandhi window are the preceding part of first predicted split $PW1$ of compound word.

The characters $NW2$ after the sandhi window are the succeeding part of second predicted split $PW2$ of compound word. Let's check the Sandhi implementation using the sequence-to-sequence model with Bi-LSTM and LSTM for the example "મનઃ + રથ = મનોરથ."

**Example: "મનઃ + રથ" = "મનોરથ"**

**Step-by-Step Explanation:**
1. **Input Sequence**:

   The input sequence is "મનઃ + રથ."This sequence is split into characters and each character is represented as a one-hot encoded vector.

2. **Encoding with Bi-LSTM**:

   The input sequence "મનઃ + રથ "is fed into the Bi-LSTM encoder layer by layer.

   Each Bi-LSTM layer processes the characters, capturing the context and dependencies in both forward and backward directions. The context captured includes how "મનઃ "& "રથ " are related and how they should combine according to Sandhi rules.

3. **End of Line Marker (<EOL>)**:

   After processing the entire input sequence, an <EOL> marker indicates the end of the input to the encoder.

4. **Decoding with LSTM**:

    The LSTM decoder takes the encoded context from the Bi-LSTM and starts generating the output sequence. The first LSTM unit receives the context and generates the first character of the output sequence. This process continues with each subsequent LSTM unit generating the next character in the output sequence based on the context and previous outputs.

5. **Output Sequence**:

    The final output sequence generated by the decoder is "મનોરથ", where the Sandhi rule has been applied to combine "મન: "& "રથ."An <EOL> marker indicates the end of the output sequence.

**Detailed Breakdown:**
**Encoder (Bi-LSTM):**

- **Step 1**: The one-hot encoded vector for "મ "is fed into the first Bi-LSTM unit.

- **Step 2**: The one-hot encoded vector for "ન "is fed into the second Bi-LSTM unit.

- **Step 3**: The one-hot encoded vector for "ઃ "is fed into the third Bi-LSTM unit.

- **Step 4**: The one-hot encoded vector for "+" is fed into the fourth Bi-LSTM unit.
- **Step 5**: The one-hot encoded vector for "ર "is fed into the fifth Bi-LSTM unit.
- **Step 6**: The one-hot encoded vector for "થ "is fed into the sixth Bi-LSTM unit.

- **Step 7**: An <EOL> marker indicates the end of the input sequence.

**Decoder (LSTM):**
- **Step 1**: The context from the Bi-LSTM encoder is feed into the first LSTM unit of the decoder.
- **Step 2**: The first LSTM unit generates the first character "મ "of the output sequence.

- **Step 3**: The second LSTM unit generates the second character "ન."

- **Step 4**: The third LSTM unit generates the third character "ો."

- **Step 5**: The fourth LSTM unit generates the fourth character "ર."

- **Step 6**: The fifth LSTM unit generates the fifth character "થ."

- **Step 7**: An <EOL> marker indicates the end of the output sequence.

**Model Execution Summary:**

    The model processes the input sequence "મન: + રથ" through the Bi-LSTM encoder to understand the context and dependencies. The LSTM decoder then generates the output sequence "મનોરથ" by applying the Sandhi rules, combining "મન:" & "રથ" into a single, grammatically correct word. This step-by-step process demonstrates how the sequence-to-sequence model effectively implements Sandhi rules in Gujarati using Bi-LSTM and LSTM layers.

**Evaluation & Outcomes:**

    Here in each stage parameter combination & results are displayed.
    (Table: Outcomes Displayed in tabular format with different comparisons)

1. **No. of Layer Based Architecture:**

| STAGE | REMARK | SR.NO. | DL MODEL | NO.OF LAYERS | VALIDATION ACCURACY | VALIDATION LOSS | Test Accuracy |
|---|---|---|---|---|---|---|---|
| 1 | NO.OF LAYER BASED ARCHITECTURE | 1 | CNN | 10 | 0.7787 | 1.7123 | 0.7129 |
| | | 2 | CNN | 20 | 0.8343 | 1.6321 | 0.8034 |
| | | 3 | RNN | 30 | 0.9574 | 1.5213 | 0.9234 |
| | | 4 | LSTM | 35 | 0.9687 | 0.8921 | 0.9545 |
| | | 5 | BI-LSTM | 40 | 0.9881 | 0.2341 | 0.9675 |

**2. Different Data Set wise:**

| STAGE | REMARK | SR.NO. | DATA SET | VALIDATION ACCURACY | VALIDATION LOSS | Test Accuracy |
|---|---|---|---|---|---|---|
| 2 | Data Set | 1 | OWN DATASET | 0.9843 | 0.6321 | 0.0803 |
| | | 2 | WORDS NET | 0.9574 | 0.5213 | 0.0924 |
| | | 3 | REAL WORLD Data Set | 0.9687 | 0.8921 | 0.9245 |

**3. Batch Size wise:**

| STAGE | REMARK | SR.NO. | MODEL | NO.OF LAYERS | VALIDATION ACCURACY | VALIDATION LOSS | Test Accuracy |
|---|---|---|---|---|---|---|---|
| 3 | Batch Size | 1 | CNN | 10 | 0.7787 | 1.7123 | 0.7129 |
| | | 2 | CNN | 20 | 0.8343 | 1.6321 | 0.8034 |
| | | 3 | RNN | 30 | 0.9574 | 1.5213 | 0.9234 |
| | | 4 | LSTM | 35 | 0.9687 | 0.8921 | 0.9545 |
| | | 5 | BI-LSTM | 40 | 0.9881 | 0.2341 | 0.9675 |

**4. No. of Epochs wise:**

| STAGE | REMARK | SR.NO. | no. of Epochs | VALIDATION ACCURACY | VALIDATION LOSS | Test Accuracy |
|---|---|---|---|---|---|---|
| 4 | NO. OF EPoch | 1 | 10 | 0.9843 | 0.2823 | 0.0803 |
| | | 2 | 25 | 0.9211 | 0.2112 | 0.9292 |
| | | 3 | 40 | 0.9375 | 0.3453 | 0.9356 |
| | | 4 | 60 | 0.9698 | 0.1235 | 0.9682 |
| | | 5 | 85 | 0.9456 | 0.0345 | 0.9344 |

The values of optimum model & Hyper parameter selected after the experiment are: Bi-LSTM model with batch size of 40, with Real world data set, epoch 60.

**Comparison with Results of different Tools of other Languages:**

Here we have compare our work with other language's implementation. The numbers shows validation Accuracy, Validation Loss & Test accuracy in below table. Actually we cannot compare our results this way because each language has different complexity in different scenario with different rule set and data. Still the best case of Gujarati sandhi is compared here in table. In this scenario the data is taken different for different languages. The main motivation of displaying this comparison is to endorse the work done in different languages' concept implementation.

| SR. NO. | NAME | VALIDATION ACCURACY | VALIDATION LOSS | TEST ACCURACY |
|---|---|---|---|---|
| 1 | JNU SANDHI TOOLS | 0.9456 | 0.2234 | 0.9343 |
| 2 | UOH SANDHI TOOLS | 0.9567 | 0.2137 | 0.9453 |
| 3 | INRIA SANDHI TOOLS | 0.9654 | 0.1234 | 0.9567 |
| 4 | Gujarati SANDHI TOOLS | 0.9789 | 0.1023 | 0.9765 |

**Conclusion:**

In this research work, we presented a model for word segmentation or "Sandhi" in Gujarati language using a purely engineering-based approach. We have tried to make a small effort to make a bridge between Gujarati Grammar's concept and implementation reality. As such the concept of Sandhi is not so well known in each South Asian language. Resources such as the "Gujarati Lexicon" and "Bhagavad-Go-Mandal" were found to be very helpful in the progression of research. The effort to design model using an extension approach for various Indian languages is expected to result in a significant lexicon resource that will be extremely valuable for natural language processing applications and machine translation. The approach proposed by the authors is capable enough to identify and perform sandhi segmentation of Gujarati language word. The approach can be

improved by improving the stemmer and lemmatization process.

**Limitation & Future Scope:**
        The accuracy after the implementation of sandhi is a huge challenge because of the language's ambiguity and law resourcefulness. Some real-time challenges are still there like some exceptional case where grammar rules are not followed properly. Each of the Gujarati language is not segment-able from available morphological rules. So, those exceptional cases must be handled properly. The windows of research are still open for researchers for all types of grammatical concept implementation and compare their results with each other using advance concepts of computer science, NLP, Machine Learning & Artificial Intelligence.

**Acknowledgment**

**References:**
1.  **Baxi**, J., & Bhatt, B. "A bidirectional LSTM-based morphological analyzer for Gujarati". *Natural Language Processing*, 1–17, (2024).

2.  **Butters**, M. Book review: The languages and linguistics of South Asia. *Lingua*.(2018). Print

3.  **Cao**, K., & Rei, M. A joint model for word embedding and word morphology. *arXiv:1606.02601,(2016).*

4.  **Yadav**, S., & Singh, R. A review of machine translation techniques for Indian languages. *International Journal of Computer Applications*, 975, 8887. (2021).

5.  **Joshi**, B. K., & Kushwah, K. K. Sandhi: The rule-based word formation in Hindi. *International Journal of Computer Science and Info. Security*,14(12), 785-790. (2016).

6.  **Kevat**, R., & Degadwala, S. A comprehensive review on Gujarati-text summarization through different features. *International Journal of Scientific Research in Computer Science, Engineering and IT*, 9(10), 301-306. (2023).

7.  **Kuncham**, P., Nelakuditi, K., Nallani, S., & Mamidi, R. Statistical sandhi splitter for agglutinative languages. In *Computational Linguistics and Intelligent Text Processing* (pp. 164–172). Springer. (2015).

8.  **Vyas**, R., & Joshi, A. An analysis of morphological features in Gujarati language processing. *International Journal of Computer Applications*, 975, 8887. (2021).

9.  **Nair**, L. R., & Peter, S. D. Development of a rule-based learning system for splitting compound words in Malayalam language. In *Proceedings of the IEEE International Conf. on Advances in Computing Communications and Informatics (ICACCI* (pp. 1-5). (2011).

10. **Yadav**, S., & Kumar, A. A study on the application of NLP in Indian languages. *International Journal of Computer Applications*, 975, 8887. (2020).

11. **Patel**, N. G., & Patel, D. B. Research review of rule-based Gujarati grammar implementation with the concepts of natural language processing (NLP). *Journal of Emerging Technologies and Innovative Research*, 5(9), 79-86. (2018).

12. **Rao**, P. S., & Kumar, A. A survey on morphological analysis of Indian

languages. *International Journal of Computer Applications*, 975, 8887. (2020).

13. **Rathi**, S., & Gupta, A. A study on the effectiveness of hybrid approaches in Gujarati text summarization. *Journal of Computer Science and Technology*, 37(1), 1-15. (2022).

14. **Sharma**, A., & Singh, R. An overview of natural language processing techniques for Indian languages. *International Journal of Computer Applications*, 975, 8887. (2019).

15. **Dave**, N. R., & Mehta, M. A. Comparative analysis of rule-based, dictionary-based,& hybrid stemmers for Gujarati language. *Comparative Analysis of Rule-Based Dicti*, 16. (2023).

16. **Singh**, A., & Kumar, R. A comparative study of stemming algorithms for Indian languages. *International Journal of Computer Applications*, 975, 8887. (2021).

17. **Modh**, J. C., & Saini, J. R. Dynamic phrase generation for detection of idioms of Gujarati language using diacritics and suffix-based rules. *International Journal of Advanced Computer Science and Applications*, 12(7), 241-247. (2021).

18. **Singh**, P., & Gupta, S. Morphological analysis of Hindi using machine learning techniques. *International Journal of Computer Applications*, 975, 8887. (2020).

19. **Soni**, H., & Patel, R. A review of idiom detection techniques in Indian languages. *International Journal of Scientific Research in Computer Science, Engineering and Information Technology*, 9(10), 301-306. (2023).

20. **Thakkar**, H., & Mehta, M. A comprehensive survey on text summarization techniques for Gujarati language. *International Journal of Computer Applications*, 975, 8887. (2022).

21. **Patel**, N., & Patel, D. Implementation approach of Indian language Gujarati grammar's concept "sandhi" using the concepts of rule-based NLP. In *Proceedings of the 2021 IEEE India Conference (INDIACom)* (pp. 1-6). (2021).

22. **Zaveri**, M., & Shah, P. A survey of sentiment analysis techniques for Gujarati language. *International Journal of Computer Applications*, 975, 8887. (2023).