# A Review On Semantic Role Labelling For Indian Languages

Rashmi R. Chouhan[1*], Dr. Charmy S. Patel[2]

1*Master Of Computer Applications Department SCET, Sarvajanik University, Surat, India, 395001 rashmi.chouhan@scet.ac.in
https://orcid.org/0000-0003-3611-8265
2Department of Computer Science SRKI, Sarvajanik University, Surat, India, 395001 charmy.patel@srki.ac.in

| ARTICLE INFO | ABSTRACT |
|---|---|
| | Semantics means the process of extracting precise meaning from text or sentences. Semantic Role Labelling (SRL) plays a crucial role in Natural Language Processing by providing insights into the underlying meaning of sentences. SRL involves assigning generic labels or roles to words within a sentence, indicating their respective semantic roles. It helps in tasks such as information extraction, question answering, educational systems, sentiment analysis, and machine translation by identifying the roles of different words in a sentence and their relationships with each other. This process enables computers to better understand and process human language, leading to more accurate and effective language understanding. SRL also powers AI-driven educational tools like intelligent tutoring systems, personalized learning, and automated assessments, making education more adaptive and effective. In this paper, we are giving a brief review of SRL system developed for different languages. This literature review focuses on key aspects of SRL like techniques used, datasets used, languages for which SRL is developed, and accuracy. Our future work is inclined towards SRL for Hindi, so our review focuses on SRL developed for Hindi, along with other Indian Languages like Urdu, Malayalam and Tamil.<br><br>**Keywords**: Natural Language Processing, Semantic Role Labelling, Semantics, Proposition Bank |

## 1. INTRODUCTION

Natural Language Processing (NLP) is an Artificial Intelligence field that mainly focuses on the interaction between computers and humans using natural language. The primary objective of NLP is to empower computers to understand, interpret, and generate human language in a way that is both meaningful and contextually relevant. One of the crucial tasks of NLP in understanding and processing human language in a more contextually aware manner is Semantic Role Labelling (SRL). SRL is a technique that identifies the predicate-argument structure of a sentence and assign the generic roles to these arguments. A variety of semantic role labels have been proposed by Fillmore's (1968, 1977) [1], widely recognized are listed in Table 1.

**Table 1.** Widely recognized Semantic Roles

| Role | Description |
|---|---|
| Agent | The doer of an action |
| Patient | The recipient of an action |
| Theme | The entity that is moved by the action |
| Experiencer | The living entity that experiences the action |
| Location | Place of an action. |
| Instrument | Intermediary used to perform an actions. |
| Source | The location/entity from which something moves. |
| Goal | The location/entity in the direction of which something moves. |

For example, "Avi gave the book to the professor", here gave is predicate (Verb), Avi is argument (Agent), book is argument (Theme) and professor is argument (Patient).

| <u>Avi</u> | <u>**Gave**</u> | <u>the book to</u> | <u>the Professor</u> |
|---|---|---|---|
| **Argument 0** | **Predicate** | **Argument 1** | **Argument 2** |
| **(agent)** | | **(theme)** | **(recipient)** |

**Figure 1.** SRL Task for English Language

| रामने | पिताजीको | पैस | दिए |
|---|---|---|---|
| **Argument 0** | **Argument 2** | **Argument 1** | **Predicate** |
| **(agent)** | **(recipient)** | **(theme)** | |

**Figure 2**. SRL Task for Hindi Language

Figure 1 and figure 2 shows how English and Hindi sentences are categorized into predicate and arguments. Further arguments can be assigned with specific roles mentioned in Table 1.

## 1.1. Applications of SRL
There are numerous real-world applications which might benefit from SRL and can improve their performance.
1.  Question Answering, [2]
SRL is the extensively studied challenge of improving predicate-argument structure for natural language words, especially verbs. In the Question Answering system the question-answer pairs are used to represent the predicate-argument structure. For example, "In 2020/21, a cricket match between India and Australia has been reshuffled, announced Monday night." The verb "announced" in this sentence would be labelled with the questions "What was announced?" and "When was something announced?" whose answers are phrases from the original sentence.

2.  Machine Translation, [3]
The task of machine translation is to change one source language word succession into another objective language word grouping which is semantically the same. SRL plays a crucial role in machine translation by providing deeper insights into the meaning and structure of sentences. By accurately identifying the roles of words and their relationships within a sentence, SRL enhances the ability of machine translation systems to produce more contextually relevant and accurate translations.

3.  Information Extraction, [4]
Information extraction is the process of extracting information from unstructured textual sources to enable finding entities as well as classifying and storing them in a database and SRL plays a crucial role in information extraction (IE) by providing a deeper understanding of the relationships between words and phrases in a sentence. By identifying the semantic roles of each constituent, SRL enables IE systems to extract more accurate and meaningful information from text.

4.  SRL in Education:
SRL can significantly enhance the educational system by improving NLP applications in various ways like Automated Grading and Assessment, Language Learning Support System, Educational Content Summarization, Natural Language Interfaces for Learning Tools, Improved Search and Query Systems and many more.

5.  Sentiment Analysis:
SRL provides valuable understanding of the roles assigned to entities and their actions in a sentence, which can help in understanding the sentiment expressed in the sentence.

6.  Other applications:
Other applications several domain like Language Learning and Education, Sentiment Analysis and Social Media Monitoring, Healthcare and Medical Records, Legal and Government Applications can be developed using SRL for Hindi Language as base application.

## 2. GENERIC STATERGY FOR SEMANTIC ROLE LABELLING

The generic strategy for text classification is depicted in Figure 3. The main steps involved after giving input are i) Pre-processing ii) Clause Boundaries Identification and Feature Selection iii) Model selection for SRL iv) SRL Output.
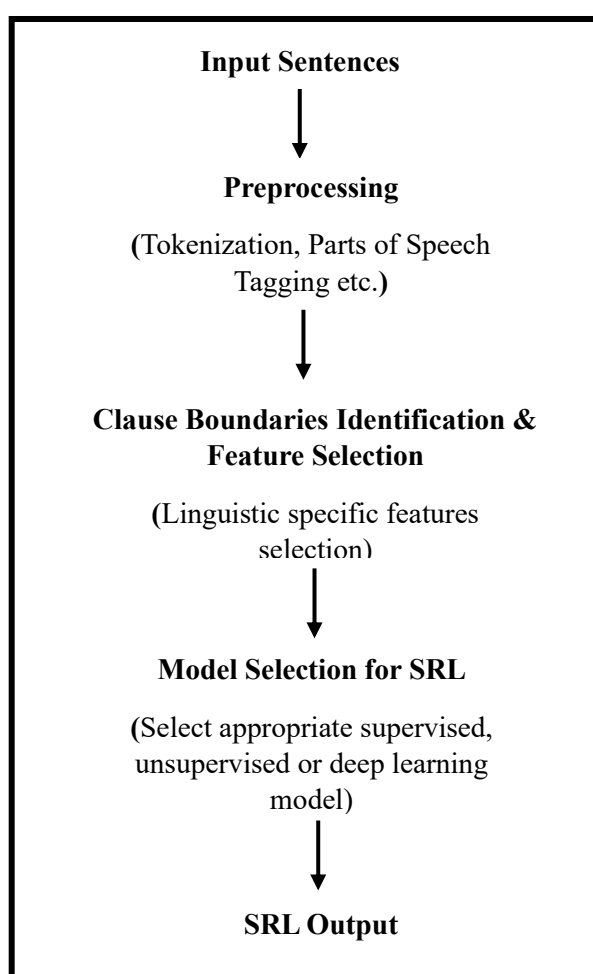Data pre-processing is a process of preparing the raw data and making it suitable for a machine learning model. It is the first and crucial step while creating a machine learning model. The pre-processing involves several tasks like tokenization, POS tagging and chunking so that sentence can become useful for SRL model. The

tokenizer divides the sentence into tokens. The Indic NLP [5] and iNLTK Library [6] provide robust tokenization support tailored for Indian languages.

Part-of-Speech (POS) tagging is essential for analysing the syntactic structure of sentences by classifying words into their respective parts of speech, such as nouns, verbs, adjectives, etc. This categorization forms the foundation for identifying main verbs and their associated roles in Semantic Role Labelling (SRL) models. Several POS taggers are available for the Indian language, including NLTK-POS-tagging [7], CRF POS Tagger 2.2 Hindi and others.

Sometimes the sentences are very large and complex. In such cases it would be better to split the complex sentences into clauses for representing the meaning. [8] The clause boundaries are identified from the clause start and clause end information. The features which can be used to find the clause start and clause end are listed below:

- Word
- POS tag
- Chunk tag
- Number of verbs in the sentence
- Number of conjunctions and subordinations in the sentence
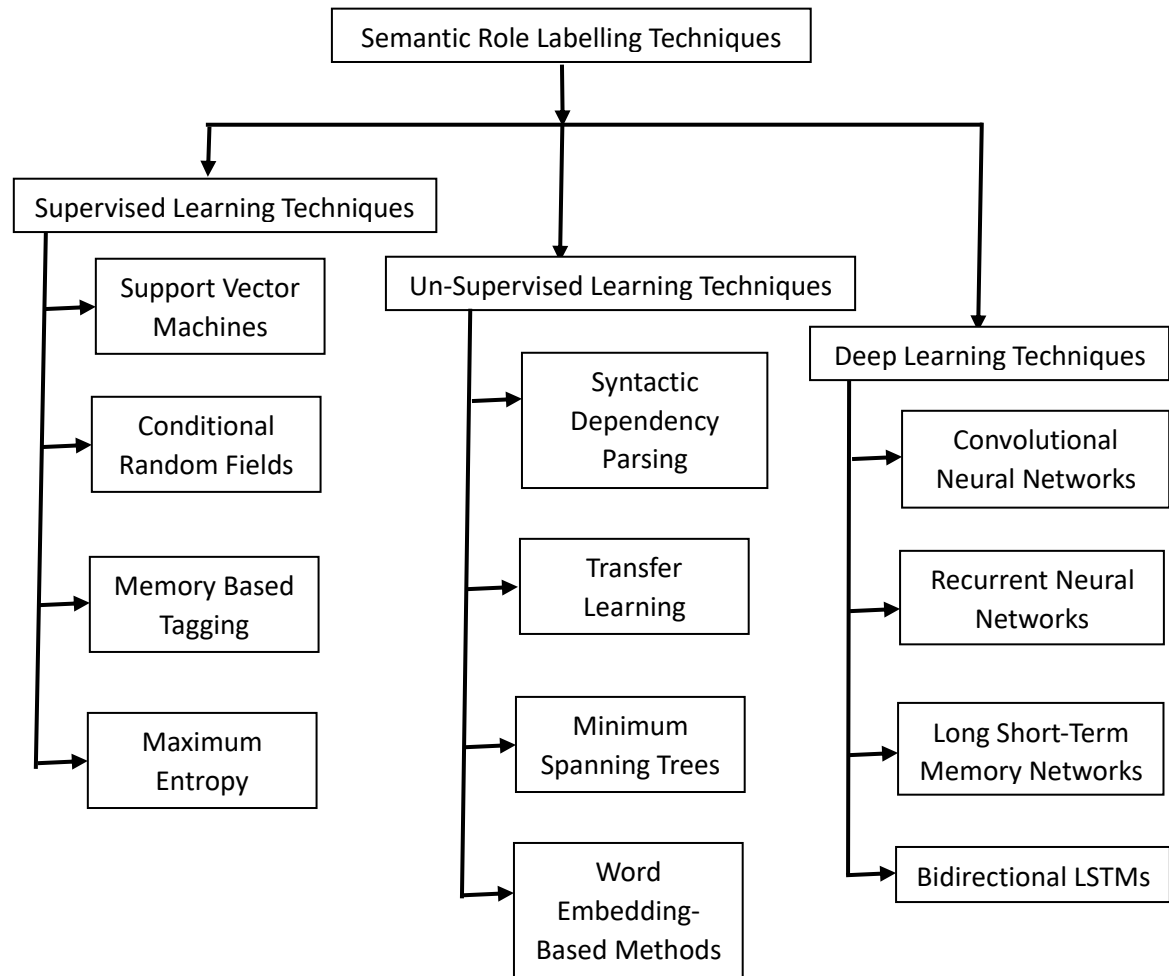- Number of negations in the sentence

**Input Sentences**

↓

**Preprocessing**

**(**Tokenization, Parts of Speech Tagging etc.**)**

↓

**Clause Boundaries Identification & Feature Selection**

**(**Linguistic specific features selection**)**

↓

**Model Selection for SRL**

**(**Select appropriate supervised, unsupervised or deep learning model**)**

↓

**SRL Output**

**Figure 3**. Generic steps for SRL

Proper set of features should be the key to SRL model, so feature selection is an important task. Features can be Predicate, Head-word, Head-word POS, Chunk Type, Dependency and many more. After all the above mentioned steps the text is ready for SRL model.

The machine learning model to be utilized in the SRL challenge needs to be carefully chosen in order to perform as accurately and optimally as possible. SRL can be built using supervised models, unsupervised models, and deep learning models. A few of the machine learning methods for SRL are depicted in Figure 4, however there may be many more. Since supervised learning models require annotated data for SRL generation, they are likely the best approach for creating SRLs for any language.

We would use Hindi for conducting our future work. Given that Hindi SRL requires annotated data for model training, supervised learning models may be the most effective approach. The primary dataset for the Hindi

SRL model will be the Hindi Proposition Bank [9], [10]. The Hindi Proposition Bank has a 400K word annotated dataset available. This dataset contains an extensive lexicon of verb frames in Hindi together with accurate annotations on the predicate-argument structure for every verb.

```
                          ┌─────────────────────────────────────────┐
                          │    Semantic Role Labelling Techniques    │
                          └─────────────────────────────────────────┘
```

Figure 4. A Few Machine Learning Techniques for SRL
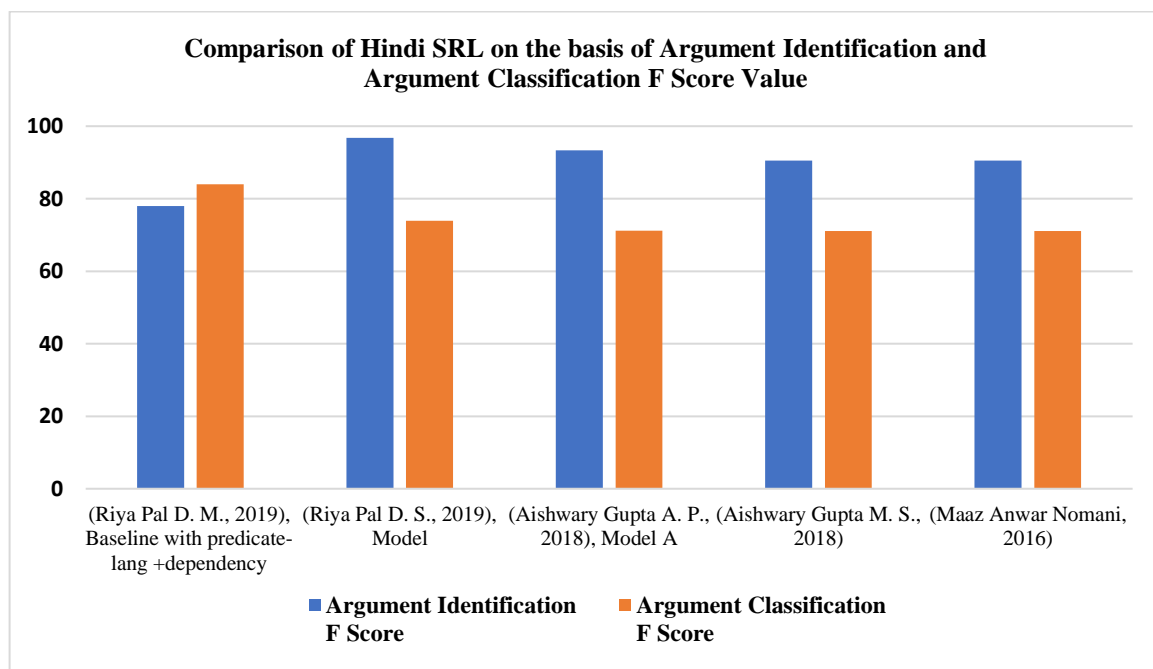
## 3. LITERATURE REVIEW

Semantic Role Labelling for Indian languages is an active area of research within the field of NLP and we can expect further advancements in NLP techniques to enhance the research in this area. Literature review on SRL for several Indian Languages is covered in this section.

### 3.1. Hindi and Urdu

a. Towards Automated Semantic Role Labelling of Hindi-English Code-Mixed Tweets [11]: The author created a technique to automate SRL for tweets with mixed Hindi and English codes. This model uses an automated labelling procedure that consists of two steps. The first step is to determine which arguments belong to which predicates in the sentence. Support vector models for binary classification are used to complete this task, making it possible to precisely detect these linguistic aspects. After arguments have been identified, the next stage is to categorize the arguments that have been found into different semantic roles. For the classification job, the model uses the Linear SVC class of Support Vector Machines (SVM), which efficiently assigns semantic roles based on the recognized parameters. Using SVM algorithms for both argument identification and role classification stages, this structured methodology guarantees accurate and automatic SRL for code-mixed tweets in Hindi-English.

b. A Dataset for Semantic Role Labelling of Hindi-English Code-Mixed Tweets [12]: In order to develop the mappings from Paninian dependency labels to Proposition Bank labels, the authors proposed a baseline rule-based framework for Semantic Role Labelling. To improve understanding, the model concentrates on syntactic differences in code-mixed data. Labelling semantic roles in Hindi-English code-mixed data had never been done before.

c. Deep learning methods for Semantic Role Labelling in Indian Languages [13]: The study shows three different models that use neural networks to increase the accuracy and productivity of SRL operations. To create embedding for dependency routes, the first model makes use of sequence modelling. This strategy allows the model to efficiently integrate the duties of argument identification and labelling within phrases, all at the same time. With the help of a bi-directional LSTM encoder, the second model may process a phrase based just on its raw tokens or words, regardless of its grammar. This architecture improves the model's independence from syntactic information in capturing contextual nuances and dependencies. The researchers add dependency labels to the final model they present in order to improve syntax awareness. This model outperforms the previous models in terms of performance by incorporating dependency labels into the bi-directional LSTM framework. Broadly speaking, the study's main goal is to use deep learning techniques to improve SRL in Indian languages. These models use neural networks to manage the intricacies contained in Indian language structures, which is a step towards more accurate and efficient SRL systems.

d. Enhancing Semantic Role Labelling in Hindi and Urdu [14]: The proposed system is supervised learning model that runs a binary classifier to classify the Arguments and then a multi-category classifier is used to classify the constituents that are labelled as arguments. In final step Support Vector Machine (SVM) Classifier is used with tuned hyper parameters for Hindi and Urdu languages. The model enhances argument identification and classification in both languages.

e. Towards Building Semantic Role Labeller for Indian Languages [15]: Proposed a supervised statistical system for identifying the semantic relationships or semantic roles for two major Indian Languages, Hindi and Urdu. Authors basically use a Logistic Regression classifier for binary classification of arguments and Linear-SVC class of SVM for multi-class strategy in argument classification. The model uses linguistic features and statistical syntactic parsing in SRL. It also uses automatic parses as features in SRL.

Figure 5 depicts the chart for comparison of Hindi SRL on the basis of Argument Identification and Argument Classification F Score Value.



**Figure 5.** Comparison of F Score value for Hindi SRL

### 3.2. Malayalam
a. Semantic role identification for Malayalam using machine learning approaches [16]: This paper focuses on semantic role identification for Malayalam language using statistical machine learning approach. Conditional Random Field (CRF) classifier is used for labelling semantic roles in Malayalam. CRF performs better as compared to Maximum Entropy, Memory Based Learning, and Support Vector Machine.

b. Semantic Role Labelling of Malayalam Web Documents in Cricket Domain [17]: Paper focuses on semantic role labelling in Malayalam cricket web documents. Memory-based language processing (MBLP) technique is used to implement the SRL approach. Model basically relied on trained POS tagger and chunker due to

lack of pre-processing tools. Model considered karta and karma karaka relations, aiming to increase accuracy.

c.  Semantic Role Labelling for Malayalam [18]: The proposed architecture consists of different pre-processing stages. After pre-processing the tagger and chunker uses a statistical model based on CRF. Model extracts meaning from sentences using Karaka theory based on Paninian Grammar and Identifies relations between nouns and verbs using case markers. The results are promising with potential for future enhancements.

d.  Semantic Parsing Approach in Malayalam for Machine Translation [19]: The author proposed algorithm for finding semantic relations using Karaka theory, as here Semantic roles are identified using Karaka theory in Paninian Grammar. It is useful for the syntax and semantic analysis of Malayalam sentences. Here Semantic parsing maps text to formal meaning representations. Shallow parsing is proposed for syntax and semantic analysis in Malayalam.

### 3.3. Tamil
a.  Entity and Verb Semantic Role Labelling for Tamil Biomedicine [20]: The author had devise a state-of-the-art SRL systems, the approach define roles to predicate terms and its constituent terms. To achieve this task MEM based classifier model is built using the features obtained from parsed input text. The MEM model is compared with Support Vector Machine and Linear Regression classifier and is found to perform better than the others.

### 3.4. Multilingual
a.  POLYGLOT: Multilingual Semantic Role Labelling with Unified Labels [21]: POLYGLOT is a multilingual SRL system for 9 languages. Focuses on multilingual semantic role labelling system with unified labels as it uses English Prop-Bank labels as universal semantic labels. Here models are trained with auto-generated Proposition Banks. The languages it includes are English, Arabic, Chinese, French, Russian, Spanish, German, Hindi and Japanese.

Table 2 shows the summary of SRL for Indian Languages. It includes the year in which study or research was done, approach, dataset on which authors had perform evaluation, Accuracy of the proposed system and limitations.

**Table 2** Semantic Role Labelling Summary

| Study | Dataset | Accuracy (%) | Limitations |
|---|---|---|---|
| **Hindi and Urdu** | | | |
| [11] | 1460 Hindi-English codemixed tweets comprising of 20,949 tokens labelled. | F1:84.00 | Limited accuracy because of typos and noisy social media data. |
| [12] | A dataset of 1460 Hindi-English code mixed tweets comprising of 20,949 tokens. | F1:84.00(Argument Identification) F1:73.93(Argument Classification) | Misclassification between different argument labels which may affects accuracy. |
| [13] | Hindi and Urdu Proposition Banks | **Hindi** F1:47.26(Model A) F1:56.55(Model B) F1:70.41(Model C) **Urdu** F1:74.88(Model A) F1: 63.15(Model B) F1: 78.07(Model C) | Model A depends on language specific feature templates. |
| [14] | Hindi and Urdu Proposition Banks | F1:90.50(Argument Identification Hindi) F1:71.12(Argument Classification Hindi) | The model does not focus on handling cases where multiple arguments of the same predicate are as-signed the same role. |
| [15] | Hindi and Urdu Proposition Banks | F1:73(Argument Identification Hindi) F1:33(Argument Classification Hindi) F1:60(Argument Identification Urdu) | Specific boundaries features had a negative impact on system results. |

| | | F1:67(Argument Classification Urdu) | |
|---|---|---|---|
| **Malayalam** | | | |
| [16] | Malayalam Sentences | F1: 60.36 | - |
| [17] | Online Malayalam Manorama News Paper cricket data. | F1: 84 | The model relies on trained POS tagger and chunker due to no standard pre-processing tools available in Malayalam that can affect the accuracy. |
| [18] | Malayalam Sentences | F1: 70 | - |
| [19] | 1000 Malayalam sentences | F1: Not Specified. Accuracy 90% | - |
| **Tamil** | | | |
| [20] | Corpus with 1180 documents | F1: 80.20 | - |
| **Multilingual** | | | |
| [21] | English, Arabic, Chinese, French, Russian, Spanish: The UN corpus of official United Nations documents.<br><br>German: The Europol corpus of European parliament proceedings and the OpenSubtitles corpus of movie subtitles.<br><br>Hindi: The Hindencorp corpus automatically gathered from web sources.<br><br>Japanese: The Tatoeba corpus of language learning. | **Arabic F Score**<br>Predicate: 93<br>Argument: 65<br>**Chinese F Score**<br>Predicate: 92<br>Argument: 82<br>**French F Score**<br>Predicate: 94<br>Argument: 80<br>**German F Score**<br>Predicate: 94<br>Argument: 81<br>**Russian F Score**<br>Predicate: 95<br>Argument: 72<br>**Spanish F Score**<br>Predicate: 95<br>Argument: 74<br>**Hindi F Score**<br>Predicate: 78<br>Argument:56 | Evaluating English Proposition Bank labels suitability for various target languages. Experimenting with constraints and heuristics to improve annotation projection quality. |

## 4. OBSERVATIONS AND FUTURE PROSPECTS

As we know, SRL is the task of identifying the predicate-argument structure in a sentence and assigning semantic roles (such as agent, patient, instrument, etc.) to the words in the sentence. This task is crucial for several natural language processing applications, as mentioned above. An extensive work is done on SRL for English, Spanish, and Chinese languages, but very little work is done for Indian languages.

As our future work will focus on SRL for the Hindi language, our observations and solutions will primarily emphasize Hindi. Below are our observations and proposed solutions for existing SRL in Hindi.
The first observation [11], [12], [13], [14], [15] from the literature review is that the techniques mainly used for Hindi SRL basically includes supervised learning models like Support Vector Machine, Linear SVM, and Logistic Regression. The second finding is that SRL for Hindi involves similar principles as in other languages, but it requires a specific model trained on Hindi text due to the linguistic differences between languages. Hindi is a postpositional language, i.e., the sentence structure of Hindi is SOV (Subject + Object + Verb). For example:

- पिताजी(S) घर (O) के अन्दर (V) हैं |
- The father (S) is inside (V) home (O).
- श्रीराम(S) सीता के बिना अयोध्या (O) कैसे जाते (V)?
- How Shri Ram(S) would have gone (V) to Adhoya (O) without Sita?
- यश (S) कार (O) चला (V) रहा है।
- Yash (S) is driving (V) the car (O).

So, the SRL model for the Hindi language requires significant modifications to the SRL model used for prepositional languages like English or we can develop a new model specifically designed for postpositional languages like Hindi.

The third observation is that the existing SRL system for Hindi [14] does not address cases where multiple arguments of the same predicate are assigned the same role. To handle such scenarios, re-ranking methods can be employed. Re-ranking can be performed using models such as feature-based, neural, context-aware, and domain-specific re-ranking models.

The fourth finding is that feature engineering methods that investigate syntactic categories, dependency roles, and head-words of chunks in the path for Hindi SRL can improve performance when paired with deep learning models [14]. But it has drawbacks of its own, such as the need for extensive processing capacity, including high-performance GPUs and enormous amounts of RAM, as well as the requirement for large volumes of annotated data to perform well.

From the overall observation of Hindi SRL, we can understand that traditional machine learning models can also give competitive performance as deep learning models. Conditional Random Fields (CRFs) and Maximum Entropy models, can be effectively used to enhance SRL performance for the postpositional languages. In particular, in the case of Malayalam SRL [16], a postpositional language, the CRF model performed better than SVM, MBL, and ME classifiers.

The findings for other Indian languages like Tamil, Malayalam, Urdu the author mostly uses supervised learning approaches for SRL.

## 5. CONCLUSION

Semantic Role Labelling (SRL) has demonstrated effective performance across numerous languages and applications like Health care, Educational, Language specific apps. However, a key challenge for most SRL methods is the necessity for annotated data, which is costly and labour-intensive to produce. This constraint remains a significant barrier to the widespread adoption of SRL across different languages. Nevertheless, several important research challenges lie ahead for advancing SRL.

In educational contexts, SRL improves language production and comprehension, fostering the growth of critical thinking and communication abilities. With the use of AI-powered teaching tools, it has the ability to completely change how students interact with texts, acquire new languages, and get feedback.

This work presents a comprehensive review of Semantic Role Labelling (SRL) for Indian languages. A range of supervised learning techniques, including Support Vector Machine (SVM), Conditional Random Field (CRF), Memory Based Tagging (MBT) and Memory Based Learning (MBL), LSTM, Bi-Directional LSTM, and Maximum Entropy (ME), have been used to construct SRL, as was indicated in section 4.

Enhancing SRL for Hindi, will be the main focus of future study. Although CRF, MBL, and ME have not yet been applied to Hindi, CRF has outperformed other supervised learning techniques for Malayalam, indicates that CRF can be a viable strategy for Hindi SRL as both are postpositional language.

## REFERENCES
1. C. J. F. a. J. B. L. C F Baker, "The Berkeley FrameNet Project," in *Proceedings of the 17th international conference on Computational linguistics-Volume 1*, 1998.
2. D. S. a. M. Lapata, "Using Semantic Roles to Improve Question Answering," in *Proceedings of the 2007 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning (EMNLP-CoNLL)*, 2007.
3. D. G. Ding Liu, "Semantic Role Features for Machine Translation," in *Proceedings of the 23rd International Conference on Computational Linguistics (Coling 2010)*, 2010.
4. M. S. S. a. O. E. Janara Christensen, "Semantic Role Labeling for Open Information Extraction," in *Proceedings of the NAACL HLT 2010 First International Workshop on Formalisms and Methodology for Learning by Reading*, 2010.
5. A. K. S. G. G. N. A. B. M. M. K. P. K. Divyanshu Kakwani, "IndicNLPSuite: Monolingual Corpora, Evaluation Benchmarks and Pre-trained Multilingual Language Models for Indian Languages," in *Findings of the Association for Computational Linguistics: EMNLP 2020*, 2020.
6. G. Arora, "iNLTK: Natural Language Toolkit for Indic Languages," in *Proceedings of Second Workshop for NLP Open Source Software (NLP-OSS)*, 2020.

7.  S. B. Edward Loper, "NLTK: the Natural Language Toolkit," in *ETMTNLP '02: Proceedings of the ACL-02 Workshop on Effective tools and methodologies for teaching natural language processing and computational linguistics - Volume 1*, 2002.
8.  A. Misha Mittal, "Clause Identification in English and Indian Languages: A Survey," *An International Journal of Engineering Sciences,* 2014.
9.  S. H. A. D. D. M. S. L. B. a. R. S. Rafiya Begum, "Dependency Annotation Scheme for Indian Languages," in *International Joint Conference on Natural Language Processing (IJCNLP)*, 2008.
10. B. N. M. P. O. R. D. S. F. X. Rajesh Bhatt, "A Multi-Representational and Multi-Layered Treebank for Hindi/Urdu," in *Proceedings of the Third Linguistic Annotation Workshop (LAW III)*, 2009.
11. D. M. S. Riya Pal, "Towards Automated Semantic Role Labelling of Hindi-English Code-Mixed Tweets," in *5th Workshop on Noisy User-generated Text (W-NUT 2019)*, Hong Kong, 2019.
12. D. S. Riya Pal, "A Dataset for Semantic Role Labelling of Hindi-English Code-Mixed Tweets," in *13th Linguistic Annotation Workshop*, Florence, Italy, 2019.
    A. P. M. S. Aishwary Gupta, "Deep Learning methods for Semantic Role Labeling in Indian Languages," in *15th International Conference on Natural Language Processing*, ICON, 2018.
13. M. S. Aishwary Gupta, "Enhancing Semantic Role Labeling in Hindi and Urdu," in *Proceedings of the LREC 2018 Workshop "The 13th Workshop on Asian Language Resources*, 2018.
14. D. M. S. Maaz Anwar Nomani, "Towards Building Semantic Role Labeler for Indian Languages," in *Tenth International Conference on Language Resources and Evaluation (LREC'16)*, Portorož, Slovenia, 2016.
15. J. S. K. T. A. Jisha P. Jayan, "Semantic role identification for Malayalam using machine learning approaches," *LOW RESOURCE MACHINE LEARNING ALGORITHMS (LR-MLA),* 2022.
16. D. A. J. A. G. SUNITHA C, "SEMANTIC ROLE LABELING OF MALAYALAM WEB DOCUMENTS IN CRICKET DOMAIN," *Journal of Theoretical and Applied Information Technology,* vol. 96, no. 8, p. 10, 2018.
17. J. S. K. Jisha P Jayan, "Semantic Role Labeling for Malayalam," in *International conference on Computational Lingustics*, 2016.
18. S. C. Shabina Bhaskar, "Semantic Parsing Approach in Malayalam for Machine Translation," *International Journal of Engineering Research & Technology (IJERT),* vol. 4, no. 7, 2015.
19. N. R. R. P. G. S. M. J. Betina Antony, "Entity and Verb Semantic Role Labelling for Tamil Biomedicine," in *7th International Conference, MIKE 2019*, Goa, Idia, 2019.
20. Y. L. Alan Akbik, "POLYGLOT: Multilingual Semantic Role Labeling with Unified Labels," in *ACL-2016 System Demonstrations*, ACL, 2016.