# Detecting and Extracting Named Entities with Particular Reference to Marathi Language

Sumalatha D. Bandari[1], Laxman L. Kumarwad[2]

[1]Department of CSE, Dr. Daulatrao Aher College of Engineering, Karad, Maharashtra, India

[2]Department of MCA, Government College of Engineering, Karad, Maharashtra, India

## ABSTRACT

The organization name, person name, location name, brand name and others are called named entities. The purpose of detecting and extracting named entities is to recognize all the named entities in the document and extracting those named entities. Detection of named entities is two step procedure- proper nouns identification and the classification of identified proper nouns. In the first step proper nouns are recognized from the text. In the second step proper nouns are classified into the different classes like the names of an organization, person, location, brand and others. Recognition of named entities is used in many applications like Natural Language Processing, Machine Translation and Machine Learning. Morphologically rich and free ordered features are present in Indian languages. Reorganization of named entities is difficult in the Indian languages like Marathi, Hindi, Urdu, Telugu and Bengali etc. The objective of this paper is to conduct the survey on recognition of named entities in different Indian languages and compared the performance metrics of different named entity approaches. Also, mentioned the challenges of Named Entity Recognition in Marathi language like morphological features, no capitalization, writing variations and ambiguity.

**Keywords :** Data Mining, J48, SMO, Naïve Bayes, Classification Algorithms

## I. INTRODUCTION

Message Understanding Conferences (MUC-6) in the year 1996 introduced the term Named Entities (NEs). The core information of the documents was present in the named entities (NEs) like organization name, person name and location names etc. In English to identify names in the document proper nouns are used. A person, place or thing names are called proper nouns. The task of detecting and extracting NEs can be described as the recognition of named entities in machine understandable form by assigning categorization tags and extraction of named entities. For example, organization names, person names and location names from corpora. Named entities detection is the most challenging task.

Named Entity Recognition (NER) is a key element in Natural Language Processing (NLP) systems for information retrieval (IR), question answering, relation extraction, etc. [1].

Language is used for human communication [2]. There are two types of communications which are verbal and non verbal. Spoken and written communications are comes under verbal communication. The structured and conventional words are used in spoken and written communication. Multilingual speakers frequently switch back and forth between languages when speaking or writing [3]. With the development of set of well defined rules inclusions and deletions are possible in the languages therefore the languages are

dynamic in nature. Indian languages are relatively free-order languages [4]. Marathi language is spoken by Maharashtra people [5]. Marathi comes under the family of Indo-Aryan languages. Marathi language is to be derived from the early forms of Sanskrit language. Marathi is the southernmost language among all the Indo-Aryan languages. Shauraseni, Magadhi and Marathi are the Prakrit languages emerged from Sanskrit. The structures of these languages are simple.

Marathi is the descendent of Maharashtra which is the Prakrit spoken by the people in the Maharashtra region. Devanagari script is used to write the languages such as Marathi, Rajasthani, Sanskrit and Nepali. The 'bALbodh' script is currently used in Marathi which is a modified version of Devanagari script. Till the eighteenth century another script called 'moDI' was used. 'moDI' script is looked similar to now a days draviDian script. In this script letters are joined together due to this feature the writing speed is more. The advantage of Dravidian script is that it is easy to read. The Devanagari script is used today which is easier to read but the disadvantage of Devanagari script is slower writing speed. Marathi script contains 52 alphabets in which 16 are vowels and 36 are consonants. Half of the words in Marathi language are derived from the Sanskrit language.

In English language, the large amount of work done for NER. The English language contains bags of resources for NER and other Natural Language Processing tasks. For the entities like organization, person and location etc., how to assign the NE tags and what type of tag is to be assigned to which type of entity is the main problem in the named entities.

There are twelve tag sets for classification of NEs.

TABLE I

INDIAN LANGUAGES NAMED ENTITY TAG SET

| Named Entity (NE) Tag | Meaning |
| --- | --- |
| NEP | Named Entity Person |
| NEL | Named Entity Location |
| NEO | Named Entity Organization |
| NED | Named Entity Designation |
| NEA | Named Entity Abbreviation |
| NEB | Named Entity Brand |
| NETP | Named Entity Title-person |
| NETO | Named Entity Title-object |
| NEN | Named Entity Number |
| NEM | Named Entity Measure |
| NETE | Named Entity Terms |
| NETI | Named Entity Time |

According to MUC structure, there are 3 types of named entities: TIMEX, NUMEX and ENAMEX [6]. Time expressions are present in TIMEX, numbers and percentages are present on NUMEX and proper names are present in ENAMEX. This experiment is interested in ENAMEX only. The proper names are categorized as follows:

Person: named person or family, for example vi. sa. khaaMdekara (वि. स. खांडेकर) or jayaMtha naaraLeekara (जयंत नारळीकर)

Organization: Corporate names, Governmental names or other organizational entities like anisa /अनिस

Location: The geologically characterized area names (urban areas, territories, districts, nations, waterways, mountains and so on.) like mahaaraaShtra/ महाराष्ट्र, bhaaratha/ भारत.

## NAMED ENTITY RECOGNITION APPROCHES

### Rule Based or Languistic Approach

The rules written manually by linguists (language experts) are used in Rule based or linguistic approach [7]. The following are the rule based NER systems

a. Gazetteer lists
b. Lexicalized grammar
c. List of trigger words

### Machine Learning Approach or Statistical Model Approach

The following are the commonly used machine learning approaches for NER

a. Conditional Random Fields (CRF)
b. Support Vector Machines (SVM)
c. Hidden Markov Models (HMM)
d. Decision Trees (DT)
e. Maximum Entropy Models (MEM)

All the machine learning approaches have their own advantages and disadvantages. The label biasing problem was not solved by Maximum Entropy Model. With the help of Markov Models Problem of sequence labelling is solved efficiently. For the development of NER system, CRF approach conditional probabilistic features and Maximum Entropy Models are very useful. Many of the related features like overlapping and non-independence are flexible to capture in CRF [8].

### Hybrid Approach

This approach is the combination of machine learning and rules based approaches. In order to improve the NER system performance, combine any two methods and make new approach which uses strongest point from each method. The hybrid approach is the combination of HMM method and Gazetteer method or HMM and MEM method.

## II. REVIEW OF LITERATURE

Vijayakrishna experimented [9] on "Domain Focused Named Entity Recognizer for Tamil Using Conditional Random Fields" in 2008. These NEs can take morphological inflection. The author builds CRF model on the noun phases of the training data. Hindi NER using the MEM given the F1 score of 71.9% for the tag set of four NE tags. General and domain specific NER was performed. To describe the named entities the finer tag set is to be needed. The system uses hierarchical tag set of 106 tags and also CRF. The troubles confronted in other machine learning systems like HMM and MEM is overcome by using CRF. Capitalization concept is not there in Tamil Language. Named Entities are represented by noun phrases. The tag sets used in this system are motivated by "Automatic Content Extraction (ACE) English Annotation Guidelines for Entities" which was developed by the Linguistic Data Consortium.

The tags used in the tag set are focused on Health Tourism domain like place, water bodies, railway stations, treatment for dieses etc. There are four levels of tag sets are present. Level-0 has 3 tags, level-1 has 22 tags, level-2 has 50 tags and level-3 has 31 tags respectively in the system. The system uses CRF++, it is an open source toolbox for straight chain CRF. The CRF model is build by the attributes extracted from the training data. To avoid ambiguity and nested tagging the system isolate the tag set into three subsets. Conditions from which the framework concludes the given phrase is the named entity but not the attributes. The conditions applied on the attributes are called the features. Noun phrase chunking is applied and only noun phrases are considered for training purpose. The attributes arrives in this system are roots of words, their Parts of Speech (POS), combined words and POS, patterns, named entity dictionary and bigram of named entity labels. Test information is processed for Morph analysis, POS and Named Person (NP) unitization.

The data to be tested is marked with everything about CRF model worked for the hierarchy of three levels. All of the three outputs measures were joined to get consolidated output. Words corpus of 94k is gathered in Tamil language for the tourism space. Transform Analysis, POS labeling, NP unitization and named substance comment are done physically on the corpus. This corpus contains concerning 20k named substances. There are 7922 Named Persons (NPs) in the test information. Completely 4059 NEs in the test information are exists. Level-1 labels are there for all of these NEs. From the 4059 NEs, 3237 NEs are having level-2 labels and 727 NEs are having level-3 labels. The purpose behind great accuracy is that labeling is done just if the premise word is to be taken from the preparation corpus. CRF is appropriate for NER in Indian dialects. Good precision is frequently got by introducing exclusively the thing phrases for each testing and training. The system obtained overall precision, recall and F-measure of 88.52%, 73.71 % and 80.44% respectively.

Asif Ekbal and Sivaji Bandyopadhyay presented [10] the "NER for Bengali & Hindi using Conditional Random Fields" in 2009. The problems of robustness and portability were exists in the rule based approach. Machine Learning approaches were effortlessly trainable and adoptable to various areas and dialects and their upkeep is less. The people with different cultures and using different languages are present in India. The web sources of name records are not accessible in Indian dialects subsequently the transliteration is used in Indian languages. By using Analysis of Variance (ANOVA) the performance is improved with the language dependent features in Bengali and Hindi languages. The expansive measure of annotated information is required to accomplish great execution for tackling NER issues. At the point when the little measure of marked information is utilized to assess the model parameters the HMMs do not function admirably. To manage different and

covering components of Indian languages MEM, CFR and SVM methods are used.

The shared task information labeled with twelve NE labels. The reason for using NE tag set is to utilize NER framework in different NLP applications and machine interpretation. To distinguish and group the maximal NEs and nested NEs the mutual undertaking was utilized. Sixteen NE labels were mapped into the four categories of NE labels, specifically Person name, Organization name, Location name and Miscellaneous. Conditional Random Fields are not the directed graphical models, the restrictively prepared probabilistic finite state automata is an uncommon instance of CRF. In noun phrase division and table extraction the CRFs had indicated achievement. For the NER problems by applying CRFs they obtained the perception arrangement is a token of sentence. By using feature induction the CRF has the choice to contain irregular components and the ability to consequently develop the most useful feature combinations. The prefixes and suffixes for all words were included in this feature.

The set of known suffixes, clue words, words, designation words and gazetteer lists are included in the language dependent features of Bengali. To recognize measurement expressions the system uses only first, middle, last names, week days and month names are used in the language dependent features of Hindi. Without the earlier information of that language, the language autonomous NE elements can be connected for NER in any language. The POS data is incorporated into the language autonomous and language subordinate components. To enhance the execution of the framework, language specific resources like lexicons, inflection lists and NER systems are used for another POS tagger. The language independent features of Bengali and Hindi are context word, word suffix and prefix features, named entity information, first word, digit features, infrequent word and part of speech information etc.

Dependent features of Bengali language have been distinguished from the Bengali news corpus. For Hindi, manually arranged the gazetteers and naturally prepared the information get from the Election Commission of India. The framework utilizes preparing set of 102,647 Bengali tokens and 452,974 Hindi tokens, advancement set of 20K Bengali tokens, 50K Hindi tokens, prepared set of 35K tokens of Bengali and 38K tokens of Hindi. The framework demonstrates the assessment of language autonomous and dialect subordinate components on the improvement sets and the evaluation results of 10-folds cross validation Test.

The system used the Bengali corpus of 122,467 tokens and the Hindi corpus of 502,974 tokens. By using twelve different NE classes these tokens are tagged. The SVM demonstrate in NER framework contains tokens 242,467 and 452,974. This is tried with the 30K and 50K tokens for Bengali and Hindi respectively. The system obtains the recall, precision, and F-score estimations of 88.7%, 80.3%, and 84.3% respectively for Bengali language and 80.5%, 74.5% and 77.4% respectively for Hindi language.

Asif Ekbal and Sivaji Bandyopadhyay presented [11] "NER system for Bengali and Hindi by using SVM model using language independent approach" in 2010. The system developed an unsupervised algorithm so as to come up with the lexical setting designs from the unlabeled news corpus of Bengali language. The lexical patterns are utilized on the grounds that the choices of SVM to support the framework execution. There are two styles of confirmations which are utilized in NER to determine the issues worried in NER. Maximum Entropy conditional models like ME Markov models and Conditional Random Fields were accounted to crush the Hidden Markov models on numerous information extraction tasks.

NER dealt with as a labeling drawback, where as each word in the sentence is allocated a mark demonstrating regardless of whether it is a piece of a named element and furthermore the entity itself.

NER systems using SVM technique is widely used by different languages and reported good accuracies. The challenges in the named entity identification of Indian languages are lack of capitalization, person names diversity, most of these names are having specific meaning in the dictionary and lengthy and difficult word forms. A little work done in Bengali NER by using pattern directed shallow parsing approach. SVM framework is additional economical than HMM or machine learning models. Gender orientation data incorporates a vital part in Hindi anyway it is not a trouble in Bengali. An unsupervised algorithm is to be utilized to create the lexical setting designs from the untagged corpus of 10 million word frames and watched the improvement in the performance.

These patterns are used in SVM based systems for post-processing the output. This system reported the result with open test and 10-fold cross validation test. To achieve the good performance stochastic model is to be used but it require large annotated corpus. HMMs do not work properly when the small amount of labeled information is used. The task data used in this system is tagged with 12 NE tags. The system considered NE tags that represent the person, organization, location names and time, number and measurement expressions. To solve the two class pattern identification problems the SVM model is to be used. High accuracy is achieved for text categorization by applying SVM. NER using SVM contains two phases: training and classification. This framework has built up an unsupervised calculation which will produce the lexical setting designs from the untagged corpus. Person, Location and Organization seed records are made. These seed lists contain 123, 87 and 32 entries. For each label embedded inside the preparation corpus, the algorithmic run creates a lexical example utilizing a setting window of six. From the 272K word forms training data 5,488 patterns are generated. Three lists are maintained for each pattern. For exactness the limit esteem is considered and the examples

underneath the edge esteems are disposed. The five percent new examples are added to the set for an every cycle of a calculation. The 10-overlap endorsement tests have demonstrated. Recall, Precision and f-score estimations of 88.69%, 80.35%, and 84.31% exclusively for Bengali and 80.48%, 74.54%, and 77.39% independently for Hindi.

B. Sasidhar presented [12] "Named entity recognition in Telugu language using language dependent features and rule based approach" in 2011. NER is a troublesome methodology in Indian dialects like Telugu, Hindi, Bengali and Urdu. The adequate gazetteers and clarified corpora are not accessible contrasted with English language. The recognizable proof of Named Substances utilizing different elements, gazetteer records utilizing language subordinate components and manage based methodologies for Telugu language. The first stage depicts the identification of nouns. The second stage distinguishes the Named Entities utilizing transliterated gazetteer records identified with distinctive Named Entity labels. NER has numerous applications in NLP viz. information arrangement, more precise web search tools, programmed ordering of reports, programmed address replying, cross language data access, and machine interpretation framework. Development of a NER framework ends up plainly difficult if legitimate resources are not accessible.

The genuine inconveniences of NER in Telugu language are lack of capitalization, agglutinative i.e. each word in Telugu language is contorted for inestimable shapes and ambiguity. The structure intertwined some gazetteer records and expansion list in the framework to increase the execution of the system. A rule based structures needs more syntactic and etymological investigation to make rules. Gazetteers are expected to separate into limited tests like suffix, prefix, setting words and so on. Gazetteers arrangement is a critical part for noun identification. NER gazette which comprises of three distinctive gazetteers Person, Location and Organization. The

transliteration depends on the phonemes and spelling. There is no quiet syllable in the composed content in these languages. Interpretation just the language changes however not the deciphered content. Transliteration system is useful for the arrangement of gazetteers records in Indian languages. Structure of Telugu language things is root stem adjacent number marker near to the case markers. Telugu language is verb last language, in each sentence last word might be a verb. Each language utilizes some particular examples which may go about as finishing words in suitable names and the rundown of this kind of words is called as postfix list. Each lingo utilizes some particular examples which may go about as suggestion words and the rundown of this kind of words is called as context records. To experiment this system window size of four is taken. Recognizable proof of root word is extremely troublesome in Telugu language. The system gathered distinctive datasets on different spaces gathered physically from different web resources, Telugu Wikipedia, Eenaadu, Andhra Prabha, Vaartha News Papers and others.

In the first stage the nouns are recognized and these distinguished nouns are given as a contribution to the second stage. Gazetteers records contain beginnings, endings, contexts and suffixes of different labels. Good performance is accomplished by this framework. 17,269 words are to be tested in that actually 16,382 NERs are present, from that 15,624 NERs are identified exactly by the system. The Named Entities identified by the system are 95.37%. Telugu language can accomplish most extreme precision by utilizing statistical learning approaches like HMM, CRF, SVM and MEM.

Nita Patil, Ajay S. Patil and B. V. Pawar present [13] "Issues and challenges in Marathi Named Entity Recognition" in 2016. Information Extraction identifies the information from the unstructured information sources. Marking NEs in linguistic communication text is critical pre-processing step helpful for information processing applications.

Implementation of Marathi NER system is difficult due to difficulties introduced in recognition. It is terribly laborious to differentiate between the characters that represent correct nouns and characters that represent traditional text. The names of the word which contain multiple meanings are called ambiguity. Word Sense Disambiguation (WSD) is needed to produce logical thinking ability to a system to see that a chunk is truly a named entity or to see the classification of a named entity. Multiple tokens are written in several ways that like abbreviations or long type, typically initial instance with descriptive long formulation followed by instances with short forms or aliases.

Parsing prediction or name constellation model is needed to predict whether or not consecutive multiple words belong to same entity. Foreign words seem in Marathi texts that are spelled in Devanagari. Recognition of foreign words is very challenging. Marathi is agglutinative language. In contrast to English prefixes and suffixes are added to root words in Marathi to create meaning contexts. Dictionaries or gazetteers contain entities without any suffix added. In Marathi suffixes are added to words in order to create the meaningful context. An elegantly composed stemmer is required for morphologically rich Marathi language to isolate the root from the addition keeping in mind the end goal to contrast the word structures and gazetteer or lexicon sections. Marathi punctuation may not be taken after totally in free-form content written work. This influences content acknowledgment frameworks. There is absence of consistency in writing spellings in Marathi. Characters composed utilizing one local character encoding may not be shown accurately by another encoding framework. On the off chance that an archive is opened on PC that does not bolster the textual style or framework utilizing which the record is composed, then content in archive is shown with disjointed characters and ends up noticeably unusable. Marathi is talked utilizing numerous lingos, for example, standard Marathi, Warhadi,

Ahirani, Dangi, Vadvali, Samavedi, Khandeshi, and Malwani in different locales of India. There are particular words used in every lingo to express the content. The fundamental procedures for named substance acknowledgment are word-level include based, rundown or word reference query and corpus based acknowledgment.

NER includes three critical subtasks such as tokenization, task of fitting tag to fragmented token and choice of right tag for a token. Vocabulary is a token with potential components and properties. Word level components portray properties of individual tokens. The techniques focuses on tokens are Regular Expressions (RE) lookup, POS tagging, Morphological analysis and shallow parsing. RE query is for the most part helpful for division of content into tokens and to portray designs. POS designs highlights word structures, for example, legitimate names that can be named substances, verbs, things and so on. Morphological examination is the way toward breaking down words into its constituents. A place of the word expression in sentence helps in understanding its part and importance in development. Shallow parsing at syntactic level can break down arrangement of words in construction. The list lookup based recognitions are Lemmatization, Stemming, Threshold edit-distance, Jaccard distance, Jarco-Wrinkle distance, Canonical normalization, Editex and Soundex.

Information driven approach is basically supervised; semi supervised or unsupervised learning techniques. HMM is a factual language display that processes the probability of an arrangement of words by utilizing a Markov chain, in which probability of next word depends on the present word. Viterbi algorithm is utilized to discover the grouping of NE classes with highest probability. HMM is regulated learning calculation. To build up a framework that can perceive named elements utilizing HMM needs tokenizer, substantial named element labeled preparing corpus, N-Gram language models, usage of

Viterbi algorithm for labeling and execution. Maximum Entropy Model registers likelihood dispersion in view of greatest entropy that fulfills the limitations set via preparing cases. Entropy is measure of vulnerability and haphazardness of the event. NER utilizing CRF depends on undirected graphical model of restrictively prepared probabilistic finite state automata. SVM is parallel characterization method used to arrange named substances. In Adaboost acknowledgment is done utilizing parallel classifiers to mark the words. Supervised learning framework utilizing decision tree utilizes data got from past, current and next word by getting some information about the history to decide conceivable yield of the model which is NE tag for a word. In bootstrapping set of seeds are utilized to begin the procedure. The framework then looks the sentences that have these seeds and tries to recognize logical signs. Clustering is unsupervised learning procedure helpful for Named Entity Recognition and Classification issues.

Important factors in unsupervised approach grouping that influences named substance acknowledgment are uncertainty in wording, contingent likelihood of the setting for a predefined semantic class and syntactic develop of the terms and setting. NER is troublesome for Indian languages and execution of Marathi language NER framework is substantially more troublesome. Different issues in Marathi NER are the intrinsic agglutinative and inflectional nature of Marathi, ambiguities in named element classes, non nearby conditions, appearances of remote words, spelling varieties and so on. F1 measure of these statistical techniques are 90.93%, 83.31%, 93.65%, 91.8%, 85%, 94.25% respectively.

Sujan Kumar Saha proposed [14] "Named Entity Recognition in Hindi using Maximum Entropy and Transliteration" in 2008. Development of NER framework is difficult if legitimate resources are not accessible. Legitimate transliteration makes the English records valuable in the NER undertakings for such languages. This framework investigated diverse elements relevant for the Hindi NER assignment and furthermore consolidated some gazetteer records in the framework to expand the execution of the framework. The authors connected the transliteration approach in Bengali NER assignment and furthermore accomplished execution change. The authors proposed a two-stage transliteration philosophy. Development of NER framework for the resource poor language is extremely difficult because of inaccessibility of legitimate resources. A portion of the resources of English language can be utilized to create NER framework for the resource poor languages. Utilization of the Indian languages in the web is practically nothing contrasted with the English language. It is conceivable to utilize these English resources if a decent transliteration framework is accessible. Transliteration is the act of translating a word or content in one written framework into another. The English names in the name records are transliterated to the intermediate alphabet.

The F-Score value accomplished by the Maximum Entropy based framework is 75.89%. At that point the transliteration based gazetteer records are fused in the framework and F-score is expanded to 81.12%. Maximum Entropy classifier is utilized to build up the framework. Maximum Entropy model has the ability to utilize diverse components to figure the contingent probabilities. In Hindi, there is no capitalization of letters to recognize formal people, location or organization names. The elements distinguished for the Hindi NER assignment are Binary Word Feature, Surrounding Words, Context Lists, Named Entity Tags of Previous Words, First Word, Containing Digit, Made up of 4 Digits, Numerical Word, Word Suffix, Word Prefix and POS Information. Maximum Entropy is an adaptable statistical model which assigns an output for each token in view of its history and components. Java based open NLP Maximum Entropy toolbox is utilized for this framework advancement. The system applied the transliteration approach for

Bengali NER task. By using this, the performance is improved in the second phase. Bilingual corpus is used to evaluate the transliteration system. This corpus contains 1,070 English-Hindi word pairs among that most of them are names. The total of 980 of them are transliterated correctly by the system and achieved an accuracy of transliteration is 91.59% for Hindi and 89.3% for Bengali.

S. Amarappa and S. V. Sathyanarayana experimented [15] on "A Hybrid approach for Named Entity Recognition, Classification and Extraction (NERCE) in Kannada Documents" in 2013. The point of this work is to build up a Hybrid model for NERCE utilizing HMM and rule-based model. The outcomes are talked about utilizing 100's of test samples. The Precision, Recall and F-measure of the system are 95.10%, 94.61%, and 94.85% respectively. Kannada is a free word order language with rich inheritance and broad accentuation. Named entity extraction in Kannada is very difficult. The Noun is identified by different elements, for example, case, number and gender. The use of named entity recognition and classification is to extract legitimate names.

Difficulties and Issues particular to Kannada language has lack of capitalization. It is Brahmi script with high phonetic trademark which could be used by NERCE framework. There is Non-accessibility of expansive gazetteer and Lack of standardization and spelling. There is number of habitually utilized words, which can likewise be utilized as names. There is lack of marked/clarified information. It is exceedingly agglutinating and bent language which requires part of artificial guidelines while separating root nouns from inflected nouns. NLP is done at various levels, for example, Phonetics and Phonology, Morphological, Syntactic Analysis, Semantic Analysis, Lexical Analysis, Discourse Integration and Pragmatic investigation. Machine Learning Algorithm HMM is utilized in this system. At that point by utilizing an arrangement of hand created rules, the perceived and ordered named elements are stemmed to extricate the root named substance. This system offers answer for some utilizations of NLP, for example, Web seeking, to check an arrangement of records written in a characteristic language and populate a database, working of helpful lexicons, building modern word processors for natural languages, Information extraction, Information recovery, Data mining, Publishing Books of Names, spots, associations and so on.

Implementation procedure of NERCE in Kannada documents using hybrid approach. Kannada editorial manager is utilized to physically make labeled Kannada corpus and spared in UTF-8 encoding group. An information base of handmade guidelines is made which comprises of root words and all possible inflections. Twelve Named Entities arranged in named entity tag sets are considered in the issue. From labeled corpus, the HMM is prepared and in the preparation arrange outflow probabilities are computed. Viterbi Algorithm is utilized to decide the Maximum probability. The state transitions for the given test output sequence is determined. The output sequence is tagged with appropriate named entity tags. Hand crafted rules are applied to extract root named entities. Assessment parameters are figured. The system is implemented using Python and Natural Language Tool Kit (NLTK) and executed on state-of-the art machine. The framework is prepared on a corpus of more than 10,000 words. The framework is tried with around 130 test samples.

Zornitsa Kozareva present [16] NER system for Spanish language using "combining different data driven systems for improving Named Entity Recognition" in 2005. This system proposed a totally programmed NER which includes distinguishing proof of legitimate names in writings and order into an arrangement of predefined classifications of interest as person names, organizations (companies, government associations, boards of trustees, and so forth.) and locations (urban areas, nations, streams, and so on). Three NE classifiers (HMM, MEM and Memory-based learner) are prepared on similar corpus information and after examination their

yields are consolidated utilizing voting system. Accuracy of 98.5% for recognition and 84.94% for classification of NE for Spanish language were accomplished. Access to learning resources in the data society is crucial to expert and self-awareness. NER has developed as a vital preprocessing tool for some NLP applications as Information Extraction, Information Retrieval and other content handling applications. Two distinctive methodologies have been created by the analysts keeping in mind the end goal to understand the NER task. The first approach depends on Machine Learning strategies; the second approach depends on Knowledge-based systems.

Jimmy L and Darvinder Kaur experimented [17] on "Named Entity Recognition in Manipuri: A Hybrid approach" in 2013. This system used hybrid approach to identify named entities in Manipuri language. The hybrid approach is the combination of statistical approach CRF and rule based approach. The supervised statistical approach for NER is used in this system. The rule based approach is used for defining the various unique word features. By using this unique word features the named entities are classified accurately by the CRF classifier. There are many challenges in Manipuri language like no capitalization, ambiguity in the named entity meaning, free order words, the class of the word is defined by using suffix, named entities are contains case markers as suffixes, highly inflectional language, complexity in stemming and limited annotated corpus, stemmer and POS tagger.

The corpus is tagged with NE tags. There are 40,000 word forms are exist in the corpus in that 10,000 are NEs. The NE tags used are PER (Person), LOC (Location), ORG (Organization), Beginning (B), Internal (I), Ending (E) multiword person, location and organization names. For labeling and segmenting of data C++ based CRF++ 0.57 package is used. The finite numbers of gazetteer lists are maintained for common location name, last names, designation etc. To create a unique identifiable feature for each multiword named entities the gazetteer lists are used.

The gazetteer list and stemming process are used to generate the training file.

The CRF features used in this system are word-[prefix/suffix], prefix, suffix, surname, first name, last name and location indicators, designation, date, currency, number in words POS tags. Among the 40,000 words of corpus 30,000 words are used in the training set. For testing purpose 10,247 words are used which contain unique words of 1,024 named entities. This system obtained the Recall, Precision and F-Score of 92.26%, 94.27% and 93.3% respectively by using hybrid approach. The label biasing problem is not present in CRF model due to this reason CRF is used in this system.

TABLE III

F-SCORE FOR DIFFERENT LANGUAGE USING DIFFERENT APPROCHES

| Language | Approach | Words | Accuracy |
|---|---|---|---|
| Tamil | CRF | 94K | 80.44% |
| Bengali | CRF | 35K | 84.3% |
| | SVM | 272K | 84.31% |
| Hindi | CRF | 50K | 77.4% |
| | SVM | 272K | 77.39% |
| | ME | | 80.00% |
| Kannada | Rule based and HMM | 10K | 94.95% |
| Manipuri | Rule based & CRF | 40K | 94.27% |
| Spanish | HMM, ME and Memory-based learner | - | 84.94% |
| Marathi | Rule based and HMM | 4,01,295 | 89.05% |

Table II shows the comparison of F-Score values of Telugu, Bengali, Hindi, Kannada and Spanish language by using CRF, SVM, MEM, machine learning techniques and hybrid methods. The observation from the above figure is the F-Score values of hybrid approaches are more as compared to individual NER approach.

This review reveals that the different types of experiments have been conducted in various languages. Some of them are using language dependent features and some are using language independent features. The gazetteers lists of the languages are required to use language dependent features. The individual NER approach cannot give good performance for detecting and extracting named entities. By using hybrid approach (combine two or more approaches) the performance of the system is improved.

## III. NER FOR MARATHI LANGUAGE

Marathi is the native and official language of the Maharashtra and spoken by the 73 million peoples according to the Census Report 2001. Marathi language has the fourth largest numbers of native speakers in India. Marathi language is speaking not only in Maharashtra but also parts of neighboring states of Goa, Dadra & Nagar Haveli, Gujarat, Daman and Diu, Chhattisgarh, Madhya Pradesh, Telangana and Karnataka [18].

Marathi is a very agglutinating and inflectional language. The most challenging sets of statistical and linguistic features are present in Marathi language. Due to this feature lengthy and difficult word forms are there in Marathi. The word begins with a root and may have a few postfixes added to one side in Marathi language. Therefore Marathi is a suffixing language. Suffixation is very difficult; it is not a simple concatenation and morphology of the language. Marathi is a resource poor language like other languages - annotated corpora, PoS taggers, good morphological analyzers, name dictionaries etc. are not so far available in the required measure. Marathi language has rich and very old literary history but the technical developments are of current origin. The web sources for name lists are accessible in English, but in Marathi language the web sources of lists are not available. Therefore, transliteration is required. NER system based on

trigram HMM model trained using preprocessed data for Marathi language [19].

Most of the research focused on resource rich language such as English language [20]. Named Entities are identified by capitalization of letters in English and also in many other languages. Upper-case, lower-case distinction is not there in Indian scripts. The capitalization feature place an important role in English language as NEs are generally capitalized in this language [21]. The capitalization concept does not exist in Indian scripts. Many of the names are common nouns. Indian names are more diverse i.e. there is lot of variation for a given named entity.

For example,

"mahaaraaShtra taaimsa/महाराष्ट्र टाइम्स" is written as

Ma Taa'/मटा etc. Developing NER systems is thus both challenging and rewarding.

The challenges in NER arise due to several factors. Some of the main factors are listed below

1. Morphologically rich - identification of root is difficult, require use of morphological analyzers.
2. No Capitalization feature - In English, capitalization is one of the main features, whereas that is not there in Marathi
3. Ambiguity - ambiguity between common and proper nouns. Example: common words such as "vinodha/विनोद" meaning Joke is a name of a person
4. Writing Variations - In the web data is that we find different people writing the same entity differently - for example: In Marathi person name - vi. sa. khaaMdekara /वि. स. खांडेकर writes like vhi. esa. khaaMdekara/ व्हि. एस. खांडेकर.

## IV. PERFORMANCE EVALUATION METRICS

The performance evaluation metrics for the data sets are precision, recall and F-Measure.

**Precision (P):** The fraction of the documents retrieved that are related to the information need for the user is called Precision.

Precision (P) = $\frac{\text{Correct Answers}}{\text{Answers Produced}}$    --------------(1)

**Recall (R):** The fraction of the documents that are related to the query that are successfully retrieved is called Recall.

Recall (R) = $\frac{\text{Correct Answers}}{\text{Total possible Correct Answers}}$    --------------- (2)

**F-Measure:** The weighted harmonic mean of precision and recall is called F-Measure. The traditional F-measure or balanced F-score is

F-Measure = $\frac{(\beta^2 + 1)\ PR}{\beta^2\ R + P}$    ------------------ (3)

β is the weighting between precision and recall. Typically the value of β=1. If β=1 the recall and precision are evenly weighted. Then the F-measure is called F1 measure. By substituting the value of β=1, the F1 - measure = 2 PR/ (P+R). There are tradeoffs between precision and recall in the performance metric.

## V. CONCLUSION

This paper discussed the advantages and disadvantages of different approaches available for NER. Rule based approach is more time consuming. Each language is having its own rules. Linguistic knowledge is required to write the rules of the languages. The NER based on rule based approach may provide high accuracy. The Machine learning approaches may not give good results because of insufficient training data. Good performance is not obtained by applying the individual NER approaches to identify the named entities. By using hybrid method that is combing the different NER approaches the performance of the system is improved. Less amount of work has been done previously in NER for Marathi language. English language NER feature like upper casing can not be utilized specifically for Marathi Language. There are many other challenges in Marathi language like ambiguity, writing variations to identify named entities.

## VI. REFERENCES

[1]  Vikas Yadav, "A Survey on Recent Advances in Named Entity Recognition from Deep Learning models", Proceedings of the 27th International Conference on Computational Linguistics, pages 2145–2158 Santa Fe, New Mexico, USA, August 20-26, 2018.

[2]  Hinal Shah, Prachi Bhandari, "Study off Named Entity Recognition For Indian Languages", International Journal of Information Sciences and Techniques (IJIST) Vol.6, No.1/2, March 2016 DOI : 10.5121/ijist.2016.6202 11

[3]  Vinay Singh, Deepanshu Vijay, Syed S. Akhtar, Manish Shrivastava, "Named Entity Recognition for Hindi-English Code-Mixed Social Media Text", Proceedings of the Seventh Named Entities Workshop, pages 27–35, Melbourne, Australia, July 20, 2018.

[4]  Shilpi Srivastava, Mukund Sanglikar & D.C Kothari. "Named Entity Recognition System for Hindi Language: A Hybrid Approach" International Journal of Computational Linguistics (IJCL), Volume (2): Issue (1): 2011.

[5]  Laxman L. Kumarwad, Rajendra D. Kumbhar and Sumalatha D. Bandari, "Present Status of Common Service Centre in Satara District of Maharashtra State (India)," 2018 8th International Conference on Cloud Computing, Data Science & Engineering (Confluence), Noida, 2018, pp. 389-393. doi: 10.1109/CONFLUENCE.2018.8442748

[6]  A. Dey, J. Abedinand & B. Purkayastha, "A Comprehensive Study of Named Entity

Recognition On Inflectional Languages", International. Journal of Advanced Research in Computer Science and Software Engineering 2014, Vol. 4, pp 696-701.

[7] Asif Ekbal et. al. "Language Independent Named Entity Recognition in Indian Languages". IJCNLP, 2008, pp 33-40.

[8] P. K. Gupta and S. Arora, "An Approach for Named Entity Recognition System for Hindi: An Experimental Study," in Proceedings of ASCNT-2009, CDAC, Noida, India, pp. 103–108.

[9] Vijayakrishna. R, " Named Entity Recognition in Tamil using Conditional Random Fields on tourism domain", Proceedings of the IJCNLP-08 Workshop on NER for South and South East Asian Languages, Hyderabad, India, January 2008. pp 59–66.

[10] Asif Ekbal and Sivaji Bandyopadhyay, "NER for Bengali & Hindi using Conditional Random Fields", LiLT Volume 2, Issue 1, November 2009.

[11] Asif Ekbal and Sivaji Bandyopadhyay "NER system for Bengali and Hindi by using SVM model", International Journal of Computer, Electrical, Automation, Control and Information Engineering Vol 4, No 3, 2010.

[12] B. Sasidhar, "Named Entity Recognition in Telugu Language using Language Dependent Features and Rule based Approach", International Journal of Computer Applications (0975 – 8887) Volume 22– No.8, May 2011.

[13] Nita Patil, Ajay S. Patil & B. V. Pawar "Issues and challenges in Marathi named entity recognition", International Journal on Natural Language Computing (IJNLC) Vol. 5, No.1, February 2016.

[14] Sujan Kumar Saha, "Named Entity Recognition in Hindi using Maximum Entropy and Transliteration", 2008

[15] S. Amarappa and S. V. Sathyanarayana "A Hybrid approach for Named Entity Recognition, Classification and Extraction (NERCE) in Kannada", Proceeding of International Conference on Multimedia Processing, Communication and Info.Tech., MPCIT, DOI: 03.AETS.2013.4.91, Association of Computer Electronics and Electrical Engineers, 2013.

[16] Zornitsa Kozareva, "NER system for Spanish language using combining different data driven systems for improving Named Entity Recognition", NLDB 2005, LNCS 3513, pp. 80–90, 2005, Springer-Verlag Berlin Heidelberg 2005.

[17] Jimmy L and Darvinder Kaur "Named Entity Recognition in Manipuri: A Hybrid approach", Gurevych, C. Biemann, and T. Zesch (Eds.): GSCL 2013, LNAI 8105, pp. 104–110, 2013.

[18] Laxman L. Kumarwad, Rajendra D. Kumbhar, "E-Governance Initiatives in Maharashtra (India): Problems and Challenges", International Journal of Information Engineering and Electronic Business (IJIEEB), Vol.8, No.5, pp.18-25, 2016. DOI: 10.5815/ijieeb.2016.05.03

[19] Nita Patil, Ajay Patil and B. V. Pawar, "HMM based Named Entity Recognition for inflectional Language", IEEE International Conference on Computer, Communications, and Electronics (COMPTELIX 2017):565-572.

[20] Kamal Sarkar, "A hidden markov model based system for entity extraction from social media english text", fire 2015. arXiv preprint arXiv:1512.03950.

[21] Sai Kiranmai Gorla, Sriharshitha Velivelli, N L Bhanu Murthy, Aruna Malapati" Named Entity Recognition for Telugu News Articles using Naïve Bayes Classifier", Proceedings of the NewsIR'18 Workshop at ECIR, Grenoble, France, 26-March-2018.