# ARPANI@BIT_DURG : KBP English Slot-filling Task Challenge

**Arpana Rawal**
Professor
Department of Computer Science & Engg.

Bhilai Institute of Technology, Durg, Chhattisgarh, India
arpana.rawal@bit
durg.ac.in

**Ani Thomas**
Professor
Department of Computer Applications

Bhilai Institute of Technology, Durg, Chhattisgarh, India
ani.thomas@bitdur
g.ac.in

**M K Kowar**
Professor
Department of Electronics & Telecommunication

Bhilai Institute of Technology, Durg, Chhattisgarh, India
mkkowar@gmail.co
m

**Sanjay Sharma**
Professor
Department of Mathematics

Bhilai Institute of Technology, Durg, Chhattisgarh, India
ssharma_bit@yah
oo.co.in

## Abstract

The present communication laid by the above mentioned TAC-SF Track participants aims to report TAC forum about the system-incorporated towards Slot-Filling task. Given a list of PER / ORG relevant attributes, the task aimed at extracting Noun-phrases that fit as assigned values to these attributes, as narrated either in one or more relevant pools of newswires or described as declarative and / or factual information in supporting knowledge-base. For accomplishing above, care was taken to identify such slot-fill patterns from every instance of relevant context whether in newswire or knowledge-base. Hence, the system relied upon generating exhaustive supporting vocabulary patterns that associate with desired slot-patterns in semantic sense. Three different kinds of sources were used to build entity-relevant vocabulary in the Task Challenge so as to search for the precise information about the entity-attributes (may be single-valued or multiple valued) mentioned in Slot-fill track task definition. The team's spirits feels elevated at the thought of using the free text from the Wikipedia pages associated with the knowledge base nodes, while building the evaluation model as it adds to the robustness and reliability of the system even at the time of inaccessibility of the Web.

## 1 Introduction

This piece of work is performed as a part of Slot-Filling task in KBP track only, that posed up the problem to design some kind of machine-assisted Question-answering system that could contextually fill an appropriate slot-value for given sets of single-valued and list-valued attributes associated with queried 'entity-name' (either PER or ORG type)given in a set of one-hundred SF-queries given at the input end. The available resources were the extended knowledge base comprising approximately 8 lakhs of the knowledge-base entity-relevant components on one side and an input corpus consisting of large number of newswires corresponding to apw, afp, xin, eng etc. news articles. For instance, source documents in slot-filling query file (available in .sgm formats) had to be fetched from TAC 2013 KBP Source Data. This also required large storage for both the collections: newspaper articles and a huge set of Wikipedia Info boxes acting as Reference knowledge bases consisting of approximately 8 lakhs of the knowledge-base entity-relevant components. The queried entities put to contextual IR task belonged to entity types: PERSON (PER) and ORGANIZATION (ORG) although the each entity type possessed a variety of info box descriptions. It was at this juncture, the very concept of entity-relevant info box was exploited

in formulating corresponding entity-type vocabularies.

## 2    ARPANI System Requirements

TAC-SF participants were initially provided with a set of 818,741 entity-descriptions in the knowledge-base (kb) corpus of the Wikipedia. The team could exploit the uniform structure of these entity descriptions irrespective of any entity type (PER / ORG / GPE / UKN), to split them into individualknowledge-base components comprising both fact-box lines and text portions.

Owing to the inherent text-mining task, the system has taken the help of computational linguistic tools like Part-Of-Speech tagger and a suitable dependency parser to explore the intra-sentential fragments. In later stages, the system has used its own designed NLP tool for generating dependency triples that precisely highlight the predominant parts viz. noun phrases, noun / verb qualifiers, verb phrases in subject and object roles.

### 2.1    Treatment Upon Slot-Filling Queries

The slot-filling queries were initially pre-processed to extract only the entity-names as the list of name-strings had three major role plays in the KBP-Slot-filling task.

### 2.1.1 Extracting pro-active knowledge-base

The participating knowledge-base component names were retrieved from SF-query statements followed by extraction of respective knowledge-base component files into a separate folder. This pre-processing step is needed further to extract distinct list of participating Info box types that lay inside the participating knowledge. Identifying such a list of info boxes shall help in generating corresponding Info box type vocabularies from the relevant knowledge-base components residing in pools of kb-corpus. The details of generating such a vocabulary are described in section 3.2.

### 2.1.2 Exploring the Participating Newswires

Secondly, these entity-mention-names could also initiate the search of the associated source document files fetched from newswire corpus. This step is necessary to extract the context in which the entity was described with reference the event narrated in the newswire. Hence, the pre-processing step to remove .xml tags and conversion of .sgm source document into .txt format was done in order to extract the participating noun-phrases that too play an inevitable role in vocabulary generation step as described in section 3.2.3.

### 2.2    Knowledge-base Refinements

Looking for the use of knowledge-base at a glance, it was realized to search for the mentioned entity-name descriptions into those kb components whose info box types matched with the info box type of the knowledge-base (kb) component in the corresponding query. However, for those SF queries with no supporting knowledge-base component, the entity-relevant description is searched in all info boxes of PER or ORG entity types; (the given queries either belong to PER / ORG entity types). For accomplishing this, two knowledge-base repositories were developed namely: PER_kbc and ORG_kbc from overall kb corpus.

With the underlying fact that both the entity type kb folders possess inclusion of numerous infobox type entity descriptions, a heuristic search was carried for extracting only those kb-components that belonged to participating infobox types. These kb-components were stored in respective folders with nomenclature format: "PER_<infobox>" or "ORG_<infobox>". For Slot-Filling (SF) queries with no supporting kb-components mentioned in their input statements, the kb-components were extracted with matches SF-names searched in '<wiki-title>' line from respective entity-type kb folders. In this way, every query has either a relevant infobox or SF-name relevant kb-folder to trigger the vocabulary generation step.

### 2.3    Fetching the Newswire contents

To extract the slot values for entity relevant attributes, the Slot_Fill names were searched in newswires, from each news category. Thus, a set of hundred SF_news folders were prepared for all the news categories. These newswires were

assimilated from all categories and are believed to possess information containment relevant to PER/ORG entity attributes.
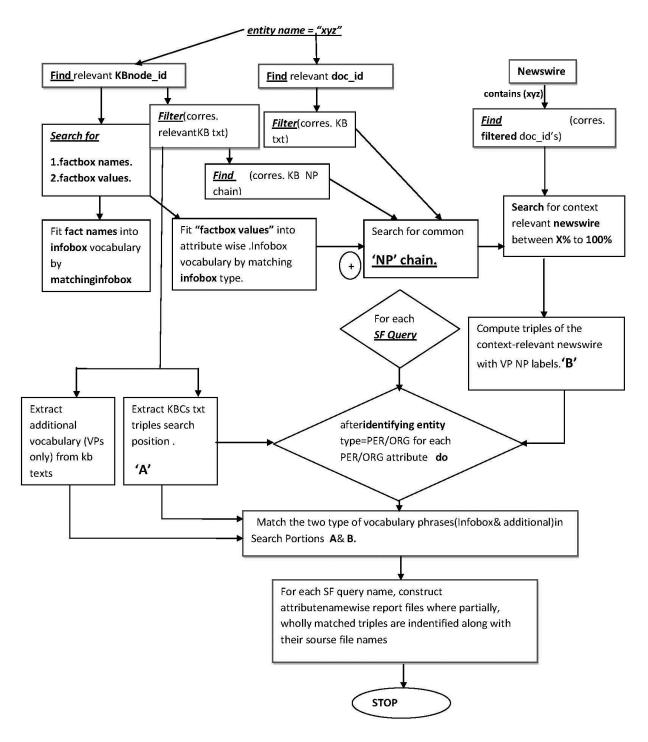
entity name = "xyz"

Find relevant KBnode_id

Find relevant doc_id

Newswire

contains (xyz)

Filter(corres. relevantKB txt)

Filter(corres. KB txt)

Find (corres. filtered doc_id's)

Search for
1.factbox names.
2.factbox values.

Find (corres. KB NP chain)

Fit fact names into infobox vocabulary by matchinginfobox

Fit "factbox values" into attribute wise .Infobox vocabulary by matching infobox type.

Search for common 'NP' chain.

+

Search for context relevant newswire between X% to 100%

For each SF Query

Compute triples of the context-relevant newswire with VP NP labels.'B'

Extract additional vocabulary (VPs only) from kb texts

Extract KBCs txt triples search position . 'A'

afteridentifying entity type=PER/ORG for each PER/ORG attribute do

Match the two type of vocabulary phrases(Infobox& additional)in Search Portions A& B.

For each SF query name, construct attributenamewise report files where partially, wholly matched triples are indentified along with their sourse file names

STOP

**Figure 1: Process Flow Diagram for TAC-KBP 2013 slot fill task challenge**

# 3    Operational Logistics

The experiments carried out by the work group outline the slot-filling task by funneling the relevant knowledge-base nodes at different conceptual depths as illustrated in *figure 1*. The experiment sequel totally rests upon the usage of computational linguistic tools for arriving at the slot-fill results which also involved the retrieval from wikipedia repository.

## 3.1    The Shallow Filtering Step

This step confines to extracting from the pools of newswire collection, only those relevant (partially / wholly) news content which portrays the involvement of the queried entity name (PER or ORG type). The extraction is implied by simply searching through the existence of desired string-pattern in the news texts. It may be noted that at this juncture, the search does not make use of context and is naïve. Similarly, in the current scenario, all possible Noun variants in the pool of Part-Of-Speech tags shall be used to extract the entity-specific noun phrases of the happenings described in newswire.

## 3.2    Vocabulary Generation Step

For the slot-filling task to accomplish the generation of supporting vocabulary strings that bind the relation between the entity names with their slot attributes (names), there was a need to understand the role of context in which the statement is constructed with the two mentioned fragments. As for any concept space to understand at machine's end the governing parameters are contributed by the extraction of Noun / verb phrases parsed from natural language text-strings. Stanford NLP group offers an efficient Part-Of-Speech Tagger among the competent ones put as open-source tools by the computational linguistic communities. The POS Tagger reads English text and assigns parts of speech to each word (and other token), such as noun, verb, adjective, etc., although generally applications use more fine-grained POS tags like 'noun-plural'. This tool is exploited to extract only noun phrases and verb phrases that describe the content relevant to the entity's attributes with the entity name itself in the underlying sentences.

### 3.2.1    Entity-Specific Vocabulary

This vocabulary is constructed by accessing to series of noun-phrase chains constructed from the participating text referred by knowledge-base id and news document id appearing in the SF-queries. The team makes use of the noun-phrase chain extraction tool that had already been developed in their TAC-KBP 2013 venture as illustrated for source document 'ENG-NG-31-142900-10121571.sgm' in *figure 2*.



**Figure 2: Entity-specific vocabulary for ENG-  NG-31-142900-10121571.sgm source document**

### 3.2.2 Fact-box Vocabulary

Some of the supporting strings could also be extracted from the set of fact-box names that were readily available from the knowledge-base components pertaining to corresponding PER / ORG entity-types that too only those kb-portions included in SF query.

### 3.2.3 Knowledge-base Vocabulary

This vocabulary is obtained by visiting each of the included simple verb or verb phrase in a knowledge-base document. When the verb (phrase) is highlighted it can be customized by extending the highlighting pointers to the left or right or both

till any word count; in this way verb phrases including all tenses / moods / auxiliaries are captured. Here human intervention is required to tap this phrase into relevant attribute vocabulary for concerned entity type using mouse clicks. For instance, if the cursor at run time traps the main verb, 'born', then capturing the neighboring fragments as 'was born on' is made to append in the attribute vocabulary of 'PER: age'. *Figure 3* shows the module screenshot while generating the supporting verb phrases forming knowledge base vocabulary for kb node "world federation of trade unions" belonging to infobox "Union".



**Figure 3: Extraction of verb-phrases for slot-attributes in knowledge-base vocabulary generation**

## 4 The Context-Filling Step

The files obtained from shallow filtering step are further examined and chosen. The noun phrases are collected from the query's supporting knowledge base as well as newswire article and they are intersected to find how many common noun phrases exist among them. These common noun phrases from knowledge base and document-id are joined to the fact box values of the knowledge base to form complete set of noun phrases. Such a chain

of noun phrases are well illustrated in *figure 4* for one of the query-associated knowledge-base nodes, SF_ENG_038. Consequently, the intersected noun phrases generated from this kb node and respective newswire article are displayed in *figure 5*.



**Figure 4: Noun-phrase generation from query-associated knowledge-base node: SF_ENG_038**



**Figure 5: Noun-phrases commonly found between knowledge-base node and news-wire document for query-id SF_ENG_038**

A matching of the complete set of noun phrases with the newswire articles obtained from the shallow filtering step allow to fetch only those news files which are having 50% or more match amongst the existing noun phrases. In this way, contextually filtered documents are retrieved from the whole set of newswires using the interface designed by the work group as shown in *figure 6*.

all these as clausal and phrasal fragments. Only those constituents appear as argument pairs that show connecting dependencies between each other. Using such combination of arguments, each fragment (paraphrased text) is mapped to its set of equivalent dependency triples. One such paraphrased text instance converted to set of triples is shown in *figure 7*.



**Figure 6: Context-relevancy test for shallow-filtered newswire documents with relevant noun-phrase chains**



**Figure 7: Extraction of vocabulary-relevant text fragments from context-filtered newswires**

## 5    The Slot-Fill Evaluation Setup

For finding the results of the runs, the contextually filtered relevant newswire documents corresponding to any query needs to be semantically analyzed and kneaded into a structure from which the slot-fillers can be found.

### 5.1 Dependency Triples As Slot-Fill Patterns

The appropriate structure chosen to represent text fragments are typed dependencies that are obtained by passing the contextually filtered files upon the Stanford Dependency Parser. As each typed dependency comprises a grammatical relation which adjoins one or more part-of-speech constituents viz. simple (common) nouns, proper nouns, noun-phrases, adjectives, pronouns, verbs, verb phrases, adverbs, prepositions, conjunctions,



**Figure 8: SF-slot value evaluation metric for TAC-KBP SF task 2013**

These in turn shall pose to help in obtaining the exact noun-phrases that enact as slot-filler values to the queried attributes for entity-names.



**Figure 9: TAC-KBP Slot-fill Results for SF query-ids 51 and 52**

**5.2 Response Generation Step**

Given the list of attributes to ignore in each SF query, the list of only those attributes is activated for which the slot-filling process is to be proceeded. The activation step triggers the creation of output files in 'edit' mode for each of such activated attribute. Hence, for 'x' no. of activated attributes in a SF query, 'x' no. of files are generated that are designed to seek the meaningful text fragments holding one or more slot-filler phrases. One can obtain such phrases by matching the vocabulary phrases from the corresponding info-box, to which the entity-name belongs. This step of computing every entity-name's info-box was an intermediary but an exhaustive computation required, as some of the queries did not have any relevant knowledge-base mentioned at the input.

For the desired attributes, each output file comprised the target sentences and the inclusive paragraph along with the metadata information i.e. filename, paragraph offset, sentence offset and the character offset information till the target sentence. Such multiple instances of target sentences are scanned for highlighting noun-phrase portions using mouse click-navigation buttons as shown in *figure 8* and the appropriate ones are picked up manually (human intervention at the final step) owing the great degree of ambiguity of English Grammar.

## 6 Conclusions

The KBP Slot-filling evaluation system encourages the research group to work at the ground level in making module design so effective that extract the newswires holding desired slot-fillers that are described well in context, on one hand and also the articles existing in the search pool having mere mention-instances of the same filler in different context, on the other hand. Both the fetched news wires should be considered as important candidates for slot-filling steps, *(see figure 9)*. It was one of the pioneering systems found that used dependency triples to filter key constituents from source text comprising noun-phrases and verb phrases and to extract the slot-fill values from them. Hence, the results were found encouraging for the moment as shown in *table 1*. The semantic structures have the scope to be used for slot-filling task in the form of dependency triples need to be revised so as to obtain higher level of system accuracy.

| ARPANI Slot-filler system Statistics | Diagnostic Response Statistics | Official Response Statistics |
|---|---|---|
| **Precision** | **0.29** | **0.28** |
| **Recall** | **0.48** | **0.5** |
| **F-Score** | **0.36** | **0.36** |

**Table 1: Slot-filling Evaluation Statistics for TAC 2013 KBP Evaluation Corpus**

research tasks are registered in Chhattisgarh Swami Vivekananda Technical University. The authors sincerely thank the aspiring team members of final year MCA students Umar, Roshan, Manoj, Sonal, Isha, and pre-final year students of BE (Information Technology) Prateek, Palash, Supriya and Nikita for their strenuous efforts in module executions to extract the TAC-KBP runs within scheduled time deadlines.

# References

*E. Agirre (University of the Basque Country)*, *A.X. Chang, D.S. Jurafsky, C.D. Manning, V.I. Spitkovsky, E. Yeh (Stanford University),* **Stanford_UBC at TAC-KBP**, Proceedings of Text Analysis Conference, 2009, organized by Information Technology Laboratory, National Institute of Standards and Technology, USA.

J. Lehmann, S. Monahan, L. Nezda, A.Jung, Y. Shi (Language Computer Corporation, USA), LCC Approaches to Knowledge Base Populationat TAC 2010, Proceedings of Text Analysis Conference, 2010, organized by Information Technology Laboratory, National Institute of Standards and Technology, USA.

H. Ji, R. Grishman, Computer Science Department Queens College and Graduate Center, New York University, Knowledge Base Population: Successful Approaches and Challenges, Proceedings of Text Analysis Conference, 2010, organized by Information Technology Laboratory, National Institute of Standards and Technology, USA.

*S. Li, S. Gao, Z. Zhang, X. Li, J. Guan, W. Xu, J. Guo (Beijing University of Posts and Telecommunications),* **PRIS at TAC 2009: Experiments in KBP Track**, Proceedings of Text Analysis Conference, 2009, organized by Information Technology Laboratory, National Institute of Standards and Technology, USA.

*V. Varma, V. Bharat, S. Kovelamudi, P. Bysani, S. GSK, K. Kumar N, K. Reddy, K. Kumar, N. Maganti (IIIT Hyderabad),* **Siel_09: IIIT Hyderabad at TAC 2009**, Proceedings of Text Analysis Conference, 2009, organized by Information Technology Laboratory, National Institute of Standards and Technology, USA.