

Query-Driven Streaming Slot-Value Extraction for Wikipedia Entities at Web-Scale

Morteza Shahriari Nia, Christan Grant, Yang Peng, Daisy Zhe Wang*, Milenko Petrovic[†]

{msnia, cgrant, ypeng, daisyw}@cise.ufl.edu, mpetrovic@ihmc.us

Abstract

In this paper we address extracting slot values from internet pertinent to entities in wikipedia, which is the most popular web-based, collaborative multilingual encyclopedia on the internet. This is comparable to Freebase in a query-driven manner. Our contributions are two fold, first we design a system to efficiently find central documents on the internet that contain directly citable content regarding a given wikipedia entity; second, we design a system to extract certain slot-values of the entity from those documents. Our results demonstrate that the system is very efficient in the web scale nature of the problem: being highly memory and I/O efficient and that we can achieve high accuracy and recall for given slot values.

Introduction

In this paper we develop an efficient query-driven slot value extraction for given Wikipedia.org (WP) entities from a timely stream of documents on internet as they are being created. Slot value extraction is formally defined as follows: match each sentence to the generic sentence structure of [Subject - Verb - Adverbial/Complement]Stevens (2008), where *Subject* represents the entity and *Verb* is the relation type we are interested in (e.g. Table 1). If a sentence matches these two components of the sentence pattern, *Adverbial/Complement* would be returned as the other side of the relation (which we refer to as slot value). Slot value extraction is a challenging task in current state of the art Knowledge Bases (KB). Popular graphical KB such as Freebase or DBpedia keep data in structured format where entities are connected via relationships (*Verb*) and the associated attributes (*Adverbial/Complement*). Our system can be used to automatically populate such KBs or even fill-in the information boxes at entity WP page itself.

Entities from three categories of *person*, *organization* and *facility* are considered. The challenges that we address are: first finding documents that contain useful information about the entity (avoid spams, documents that do not have direct information about the entity even

though they mention them, etc), the scale where *stream processing* nature of the system makes it suitable to avoid batch processing natures of Hadoop and operate in the realm of streams such as twitter storm¹ which similarly processes unbounded streams of twitter data or Spark Streaming². As other Natural Language Processing tasks, precision and recall of the extent we can extract slot values are very important metrics that we will discuss later on. Having to balance between infamous and obscure entities, dealing with slots that can be very broad (e.g. things that an entity is affiliated with) or very specific (such as cause of death of a *person* entity).

To take on an example assume we are analyzing a sentence from the internet such as "Boris Berezovsky made his fortune in Russia in the 1990s when the country went through privatisation of state property and 'robber capitalism', and passed away March 2013.". First we have to pay attention that there are two *Boris Berezovsky* entities in WP, one a businessman and the other a pianist. Any slot value extraction shall take this into account and try to come up with a viable distinguishing policy (Entity Resolution). Then, we match the sentence to find a slot value such as *DateOfDeath* valued at March 2013. Other examples that this system can answer could be 'Who a person has met during a certain period of time?', 'Who are the employees of this organization?', 'Who has met who in this facility at a certain time?'.

An important challenge in maintaining WP, the most popular web-based, collaborative, multilingual KB on the internet, is making sure its contents are up-to-date. Presently, there is considerable time lag between the publication date of cited news and the date of an edit to WP creating the citation. The median time lag for a sample of about 60K web pages cited by WP articles in the *living-people* category is over a year and the distribution has a long and heavy tail Frank et al. (2013). Also, the majority of WP entities have updates on their associated article much less frequently than their mention frequency. Such stale entries are the norm in any

¹<http://storm-project.net/>

²<http://spark.incubator.apache.org/>

Table 1: Ontology of Slots

PER	FAC	ORG
Affiliate AssociateOf Contact_Meet_PlaceTime AwardsWon DateOfDeath CauseOfDeath Titles FounderOf EmployeeOf	Affiliate Contact_Meet_Entity	Affiliate TopMember FoundedBy

large KB because the number of humans maintaining the KB is far fewer than the number of entities. Reducing latency keeps KBs and WP relevant and helpful to its users. Given an entity page, possible citations may come from a variety of sources. The actual citable information is a small percentage of the total documents that appear on the web. We develop a system to read streaming data and filter out articles that are candidates for citations.

Our system contains three main components. First, we pre-process the data and build models representing the KB entries. Next, we use the models to filter a stream of documents so they only contain candidate citations. Lastly, we process sentences from candidate extractions and return slot values. Overall, we contribute the following:

- Introduce a method to build models of name variations
- Build a system to filter a large amount of diverse documents
- Extract, infer and filter entity-slot-value triples of information to be added to KB

System

In this section, we introduce the main components of the system. Our system is built with a pipeline architecture in mind giving it the advantage to run each section separately to allow stream processing without blocking the data flow of components (Figure 1). The three logical components include sections on *Model* for entity resolution purposes, *Wiki Citation* to annotate cite-worthy documents, *Slot Filling* to generate the actual slot values.

To walk you through the steps we take, assume we only care about on single WP entity, the first step is to extract aliases of the entity. We use several approaches to get as many viable aliases as possible. Then we look into the stream of content that is being generated on the internet, apply two levels of filtering to finally end up with the documents are central to that entity. To extract the relevant slot values we perform pattern matching in each sentence or coreferent sentence to see if we can find a match. As a match is found from the content of the sentence to the patterns that we have generated regarding slot name, the associated slot value is extracted

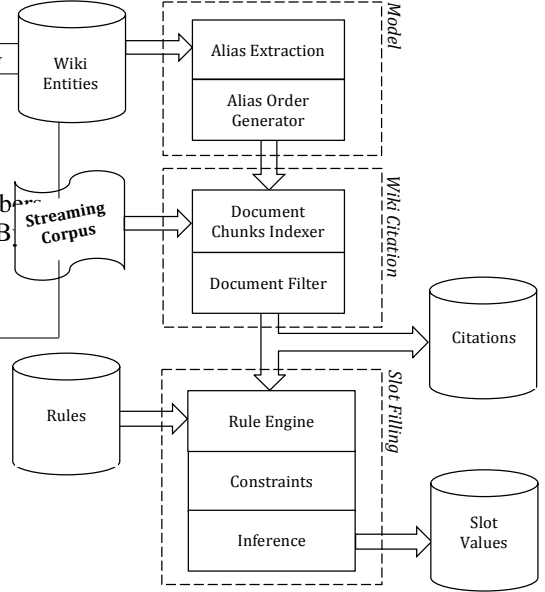


Figure 1: System Architecture. Components are logical groups noted with dotted boxes.

as a final result.

Model

We use Wikipedia API to get these aliases automatically. This is done by retrieving backlink references (redirects of a wiki entity), e.g. William Henry Gates is an alias for Bill Gates in WP as a backlink reference. Unfortunately this is not good enough and to enhance recall we need more aliases. To have better use of a wiki page we parse HTML DOM of the page, then use regular expressions to extract the bold phrases of the first paragraph as alias of the actual entity. Based on our observation this is a very accurate heuristic and provides us with lots of famous aliases of the entities. As an example of when this might not work, is that there might be occasions that some other topic is written in bold typesetting in the first paragraph apart from the entity aliases itself but these are very rare.

Once aliases are available we pass them through rules of generating proper name orders which will produce various forms of writing a name. As a basic example Bill Gates can be written as Gates, Bill also. This will allow the system to capture various notation forms of aliases. We refer to this part as *Alias Order Generator*.

Cumulative Citation Recommendation. The main goal of CCR is to have an aggregate list of documents that are worthy of being cited in a Wikipedia page. We perform exact string matching and treat all the documents that mention an entity equally likely to be citable. One of the reasons for this is that in former TREC KBA reports Frank et al. (2013) there were observations of how non-mentioning documents have a low chance of being citable in Wikipedia. So we take on that and ignore non-citing documents.

Streaming Slot Filling. The purpose of SSF is to ex-

tract proper values for relations of interest, which can be found in Table 1. This is called Stream Slot Filling because data is being generated as time goes on and for each extraction we should only consider current or past data. In Figure 1 we refer to this as *Streaming Slot Value Extraction*. Stream slot filling is done by pattern matching documents with manually produced patterns for slots of interest. The way we do this is by observing a sentence that has a mention of the entity or one of its coreferences. An anchor word in the sentence related to the slot name is located and we match either left or right of the anchor word for potential slot values.

Post Processing Algorithm. The SSF output of many extractions is noisy. The data contains duplicates and incorrect extractions. We can define rules to sanitize the output only using the information present in the SSF file. The file is processed in time order, in a tuple-at-a-time fashion to minimize the impact on accuracy. We define two classes of rules deduplication rules and inference rules. In our diagram we refer to this component as **High Accuracy Filter**.

Implementation

We extract aliases for entities from Wikipedia automatically both using API and using the actual page content, then apply pattern matching rules for slot value extraction. Our contribution is that we perform pattern matching that conforms to each slot value along with post-processings to eliminate noisy outputs.

Alias Generation

Cumulative Citation Recommendation

Our pipeline of processing the corpus consists of a two layer indexing system referred to as *Chunk Files Index Generator* and *StreamItems Index Generator*. Chunk Files Index Generator will generate indexes of the chunk files that contain a mention of any of the desired entities. StreamItems Index Generator on the other hand will index StreamItems that contain a mention of a given entity respectively. This two level indexing will eliminate the need to process each and every Chunk-File/StreamItem for every entity. The reason for splitting this task into two steps is that not all chunk files contain any mention of the entities and we want to get rid of them as soon as possible. Chunk Files Index Generator which discards non-mentioning chunk files and will stop further processing a chunk file as soon as it finds a mention there. Each chunk file can contain up to thousands of SIs which can be so time consuming if we were to process them in our Java base code. Processing StreamItems on the other hand is done in Java with ideas in mind for later on extensibility by adding other Java libraries.

Streaming Slot Value Extraction

In the data set, we are given a date range of documents as training data. Instead of building a classifier we use pattern matching methods to find corresponding

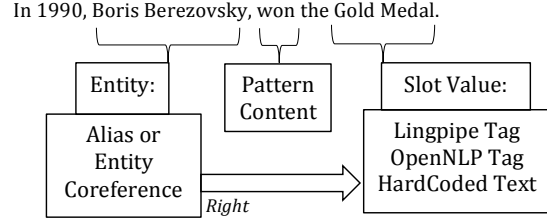


Figure 2: Pattern Matching with Slot Value on the Right Side of Entity.

slot values for entities. Pattern matching is simple to manipulate results and implement. Additionally, a classifier approach is more difficult to evaluate and explain results.

With the StreamItem indexes generated by the CCR, we first fetch the sentences containing entites by using alias names and coreference information provided by Lingpipe tags. Then use these senteces to match patterns and when patterns matched, generate SSF results.

Algorithm 1 Streaming Slot Value Extraction Pseudocode

List of entities $\mathcal{E} = \{e_0, \dots, e_{170}\}$

List of patterns $P = \{p_0, \dots, p_{|P|}\}$

List of streamitems containing entities
 $\mathcal{S} = \{s_0, \dots, s_{|S|}\}$

```

for  $si \in \mathcal{S}$  do
  for  $sentence \in si$  do
    for  $entity \in \mathcal{E}$  do
      if Contains( $sentence, entity$ ) then
        for  $pattern \in P$  suitable for  $entity$  do
          if Satisfies( $sentence, pattern$ ) then
            Emit( $sentence, pattern$ )

```

Format of patterns A pattern is defined as a record representing knowledge going to be added to a knowledge base. A pattern P is represented as a five-tuple $P = \langle p_1, p_2, p_3, p_4, p_5 \rangle$.

The first value, p_1 represents the type of entity. These entity types are in the set $\{FAC, ORG, PER\}$ where FAC represents a type of facility, ORG represents an organization and PER represents a person. FAC, ORG and PER are Lingpipe entity types. The p_2 represents a slot name. A list of slot names is present in Table 1. The third element p_3 is the pattern content. This is a string found in the sentence. The extractor looks for this exact string or pattern in a sentence. The pattern evaluator uses a direction (left or right) found in p_4 to explore sentence. The final element p_5 represent the slot value of a pattern. This The type of slot value may be the entity type tagged by Lingpipe, a noun phrase (NP) tagged by OpenNLP or a hard-coded phrase. For these three kinds of patterns, we implement them in different ways accordingly. Next, we explain the patterns with more details, an example can be found in Figure 2.

Types of patterns There are three types of patterns distinguished by different types of slot values in the patterns. The matching methods using these three types of patterns are implemented according to the different information and structures of slot values.

Type I. This pattern type is driven by the slot value type, a pattern tagged by Lingpipe. For example, pattern $\langle \text{PER}, \text{FounderOf}, \text{founder}, \text{right}, \text{ORG} \rangle$. PER means that the entity we are finding slot values for a PER entity; FounderOf means this is a pattern for FounderOf slot. *founder* is the anchor word we are match in a sentence; *right* means that we are going to the right part of the sentence to match the pattern and find the slot value; ORG means the slot value should be a ORG entity.

Type II. This pattern type is unique because it only looks for a slot value tagged as noun phrase (NP) by OpenNLP. For example, pattern $\langle \text{PER}, \text{AwardsWon}, \text{awarded}, \text{right}, \text{NP} \rangle$. This pattern can be interpreted as that we are looking for a noun phrase after the *awarded* since that noun phrase may represent an award. Titles and awards are usually not the Lingpipe entities, hence the use of the OpenNLP noun phrase chunker to fetch the noun phrases.

Type III. Some relations are best discovered by hard coding the slot values. Examples of these include time phrases: $\langle \text{PER}, \text{DateOfDeath}, \text{died}, \text{right}, \text{last night} \rangle$. In this pattern, *last night* means we are looking for exactly the phrase *last night* to the right of *died*. This pattern is inspired by the intuition that in news articles, people often mention that somebody died last night instead of mentioning the accurate date information and Lingpipe tends not to tag phrases like *last night* as a DATE entity.

A short discussion We sampled documents from the training data period to generate an initial set of patterns. We then use these patterns to generate SSF results. By manually looking at these results, we prune some patterns with poor performance and add more patterns that we identified from these results. We use several iterations to find the best patterns. We found that it is time consuming to identify quality pattern.

We found three major classes of accuracy errors: incorrect entities selected, incorrect tags by Lingpipe and incorrect pattern extractions. The first issue is ameliorated by generating better aliases (Section). And we use post-processing to reduce the second and third types of errors (Section). We didn't use more advanced NLP packages such as Stanford NLP because of the large size of the data set. The post-processing step to improve the results is discussed in the next section.

High Accuracy Filter

The SSF output of streaming slot value extraction is noisy. The data contains duplicates and incorrect extractions. We can define rules to sanitize the output only using the information present in the SSF file. The file is processed in time order, in a tuple-at-a-time fashion to

minimize the impact on accuracy. We define two classes of rules: *deduplication* and *inference* rules.

The output contains many duplicate entries. As we read the list of extracted slots we create rules to define "duplicate". Duplicates can be present in a window of rows; we use a window size of 2 meaning we only be adjacent rows. Two rows are duplicates if they have the same exact extraction, or if the rows have the same slot name and a similar slot value or if the extracted sentence for a particular slot types come from the same sentence.

New slots can be deduced from existing slots by defining inference rules. For example, two slots for the task are "FounderOf" and "FoundedBy". A safe assumption is these slot names are biconditional logical connectives with the entities and slot values. Therefore, we can express a rule "X FounderOf Y" equals "Y FoundedBy X" where X and Y are single unique entities. Additionally, we found that the slot names "Contact_Meet_PlaceTime" could be inferred as "Contact_Meet_Entity" if the Entity was a FAC and the extracted sentence contained an additional ORG/FAC tag. We also remove erroneous slots that have extractions that are several pages in length or too small. Errors of extracting long sentences can typically be attributed to poor sentence parsing of web documents. We have some valid "small" extractions. For example a comma may separate a name and a title (e.g. "John, Professor at MIT"). But such extraction rules can be particularly noisy, so we check to see if the extracted values have good entity values.

Evaluation

Our system was developed on a 32-core server described in Table 3. The corpus is a snapshot of the web in English. Each document is annotated using lingpipe and is called StreamItem, a bundle of StreamItems are put together and serialized as Apache Thrift objects, then compressed using xz compression with LempelZivMarkov chain algorithm (LZMA2) and finally encrypted using GNU Privacy Guard (GPG) with RSA asymmetric keys. The total size of the data after XZ compression and GPG encryption is 4.5TB and just over 500M StreamItems s3. Data is stored in directories the naming of which is date-hour combination: from 2011-10-05-00 (5th of October 2011, 12am) until 2013-02-13-23 (13th of February 2013, 11pm), which consists of 11952 date-hour combinations, 4 months of data (October 2011 - February 2012) is for training purposes. This corpora consists of various media types the distribution of which can be found in Table 2. To have a sense of the scale of objects and compression as an example a 6mb gpg.xz files would become 45 mb thrift objects which can contain a couple of thousand StreamItems depending on their size. Some of the documents have null values for their annotation fields. The first portion of the data which ranges from October 2011 to February 2012 is considered as training data. The source code of our system is stored as an open source project where

Table 2: Document Chunks Distribution

# of Documents	Document Type
10988	arxiv (full text, abstracts in StreamItem.other_content)
34887	CLASSIFIED (spinn3r)
77674	FORUM (spinn3r)
12947	linking (reprocessed from kba-stream-corpus-2012, same stream_id)
141936	MAINSTREAM_NEWS (spinn3r)
4137	MEMETRACKER (spinn3r)
280629	news (reprocessed from kba-stream-corpus-2012, same stream_id)
6347	REVIEW (spinn3r)
688848	social (reprocessed from kba-stream-corpus-2012 plus extension, same stream_id)
740987	WEBLOG (spinn3r)

enthusiasts can also contribute to git, also the relevant discussion mailing list is accessible here goo.

We have 135 extraction patterns coverin each slot-name/entity-type combinations. Our final submission was named *submission_infer*. Our results are as follows: Document extraction using query entity matching with aliases, sentence extraction using alias matching and co-reference. Slot extraction using patterns, NER tags and NP tags. 158,052 documents with query entities, 17885 unique extracted slot values for 8 slots and 139 entities, 4 slots and 31 entities missing.

On the performance of our initial submission run we performed random sampling via two processes, the results of which are according to Table 4. You can view that we have had an accuracy of around 55%, and about 15% wrong entity identified and 30% incorrect value extracted across all entities and slot types. Most of our issues for this submission were regarding poor slot value extraction patterns and incomplete aliases which were tried to be mtigated later on. For our final submission, we provide a more detailed statistics, which has been elaborated in Table 5 and Table 6. Table 5 shows the extent of search outreach for each slot name. You can see that *Affiliate* has been the slot name with highest hits and *CauseOfDeath* our lowest hit with 0 instances found matching our patterns, after that *AwardsWon* has been the next with 38 instances found. Affiliate is a very generic term and extracting real affiliates can be quite challenging using the extraction patterns provided. This can lead to noisy results. On the other hand for more precise terms our accuracy increases but we have less recall. Table 6 addresses the relative accuracy measure per slot value. There you can view that we have had the highest accuracy of 63.6% for *AssociateOf* and the lowest of 1% - 5% for *Affiliate*, *Contact_Meet_PlaceTime* and *EmployeeOf*.

Discussions & Conclusion

Table 5 show a varied distribution of extracted slot names. Some slots naturally have more results than other slots. For example, *AssociateOf* and *Affiliate* have more slot values than *DateOfDeath* and *CauseOfDeath*, since there are only so few entities that are deceased. Also, some patterns are more general causing more ex-

Table 3: Benchmark Server Specifications

Spec	Details
Model	Dell xxx 32 cores
OS	CentOS release 6.4 Final
Software Stack	GCC version 4.4.7, Java 1.7.0_25, Scala 2.9.2, SBT 0.12.3
RAM	64GB
Drives	2x2.7TB disks, 6Gbps, 7200RPM

Table 4: SSF Performance Measure on *initial_submission*

	Correct	Incorrect Entity name	Incorrect Value
Sampling #1	55%	17%	27%
Sampling #2	54%	15%	31%

tractions. For example, for *Affiliate*, we use *and*, *with* as anchor words. These words are more common than *dead* or *died* or *founded* in other patterns.

When we evaluate the results of slot extraction, we find three kinds of problems for accuracy: 1) wrong entities found; 2) wrong tags by the Lingpipe; 3) wrong results matched by the patterns. We also have recall problems: 1) not enough good alias names to find all the entities. 2) not enough and powerful patterns to capture all the slot values. We will use entity resolution methods and other advanced methods to improve the accuracy and recall of entity extraction part.

For slot extraction, to improve the performance, we need: 1) Using multi-class classifiers instead of pattern matching method to extract slot values in order to increase both recall and accuracy for slots “*Affiliate*”, “*AssociateOf*”, “*FounderOf*”, “*EmployeeOf*”, “*FoundedBy*”, “*TopMembers*”, “*Contact_Meet_Entity*” and so on. 2) For special slots, like “*Titles*”, “*DateOfDeath*”, “*CauseOfDeath*”, “*AwardsWon*”, using different kind of advanced methods, e.g. classifiers, matching methods. 3) Using other NLP tools or using classifiers to overcome the drawbacks of the LingPipes inaccurate tags. The first and second tasks are the most important tasks we need to do.

Table 5: Recall Measure on *submission_infer*: Generic slot names like affiliate had the most recall, compared to less popular slot names e.g. DateOfDeath

Slot Name	Total instances of slot value found	# of entities covered by slot value
Affiliate	108598	80
AssociateOf	25278	106
AwardsWon	38	14
CauseOfDeath	0	0
Contact_Meet_Entity	191	8
Contact_Meet_PlaceTime	5974	109
DateOfDeath	87	14
EmployeeOf	75	16
FoundedBy	326	30
FounderOf	302	29
Titles	26823	118
TopMembers	314	26

Table 6: SSF Accuracy Measure on *submission_infer*: Accuracy of AffiliateOf was the best and Affiliate applied poorly due to ambiguity of being an affiliate of somebody/something

Slot Name	Correct	Incorrect Entity name	Incorrect Value
Affiliate	1%	95%	5%
AssociateOf	63.6%	9.1%	27.3%
AwardsWon	10%	10%	80%
CauseOfDeath	0%	0%	0%
Contact_Meet_Entity	21%	42%	37%
Contact_Meet_PlaceTime	5%	20%	85%
DateOfDeath	29.6%	71%	25%
EmployeeOf	5%	30%	65%
FoundedBy	62%	17%	21%
FounderOf	50%	0%	50%
Titles	55%	0%	45%
TopMembers	33%	17%	50%

About 50% of twitter entities are not found by the system. One reason is those entities are not popular. For example, a 'Brenda Weiler' Google search result has 860,000 documents over the whole web. For our small portion of the web it might make sense. The histogram of the entities shows that more than half of the entities have appeared in less than 10 StreamItems. A good portion have appeared only once.

Alltogether, We experimented through different tools and approaches to best process the massive amounts of data on the platform that we had available to us. We generate aliases for wikipedia entities using Wiki API and extract some aliases from wikipedia pages text itself. On twitter entities we extract aliases manually as it is part of the rule of the KBA track. We process documents that mention entities for slot value extraction. Slot values are determined using pattern matching over coreferences of entities in sentences. Finally post processing will filter, cleanup and infers some new slot values to enhance recall and accuracy.

We noticed that some tools that claim to be performant for using the hardware capabilities at hand sometimes don't really work as claimed and you should not always rely on one without a thorough A/B testing of performance which we ended up in generating our in-house system for processing the corpus and generating the index. Furthermore, on extracting slot values, pattern matching might not be the best options but definitely can produce some good results at hand. We have plans on generating classifiers for slot value extraction purposes. Entity resolution on the other hand was a topic we spent sometime on but could not get to stable grounds for it. Entity resolution will distinguish between entities of the same name but different contexts. Further improvements on this component of the system are required.

Acknowledgements

Christan Grant is funded by a National Science Foundation Graduate Research Fellowship under Grant No. DGE-0802270.

References

Ace (automatic content extraction) english annotation guidelines for events. http://projects.ldc.upenn.edu/ace/docs/English-Events-Guidelines_v5.4.3.pdf.

Frank, J. R.; Kleiman-Weiner, M.; Roberts, D. A.; Niu, F.; Zhang, C.; Ré, C.; and Soboroff, I. 2013. Building an entity-centric stream filtering test collection for trec 2012. In *21th Text REtrieval Conference (TREC'12)*. National Institute of Standards and Technology.

Gatordsr opensource project. <https://github.com/cegme/gatordsr>.

Gatordsr mailing list. <https://groups.google.com/forum/#!forum/gatordsr>.

Ji, H., and Grishman, R. 2011. Knowledge base population: Successful approaches and challenges. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies - Volume 1*, HLT '11, 1148–1158. Stroudsburg, PA, USA: Association for Computational Linguistics.

Trec kba stream corpora. <http://aws-publicdatasets.s3.amazonaws.com/trec/kba/index.html>.

Stevens, S. 2008. Introduction to sentence patterns. In *Tutoring and testing at UHV*. University of Houston - Victoria.

Tac kbp slots. http://www.nist.gov/tac/2012/KBP/task_guidelines/TAC_KBP_slots.html.