

SlugNERDS: A Named Entity Recognition Tool for Open Domain dialogue Systems

Kevin K. Bowden, Jiaqi Wu, Shereen Oraby, Amita Misra, and Marilyn Walker

Natural Language and dialogue Systems Lab

University of California, Santa Cruz

{kkbowden, jwu64, soraby, amisra2, mawalker}@ucsc.edu

Abstract

In dialogue systems, the tasks of named entity recognition (NER) and named entity linking (NEL) are vital preprocessing steps for understanding user intent, especially in open domain interaction where we cannot rely on domain-specific inference. UCSC's effort as one of the funded teams in the 2017 Amazon Alexa Prize Contest has yielded Slugbot, an open domain social bot, aimed at casual conversation. We discovered several challenges specifically associated with both NER and NEL when building Slugbot, such as that the NE labels are too coarse-grained or the entity types are not linked to a useful ontology. Moreover, we have discovered that traditional approaches do not perform well in our context: even systems designed to operate on tweets or other social media data do not work well in dialogue systems. In this paper, we introduce Slugbot's Named Entity Recognition for dialogue Systems (SlugNERDS), a NER and NEL tool which is optimized to address these issues. We describe two new resources that we are building as part of this work: SlugEntityDB and SchemaActuator. We believe these resources will be useful for the research community.

Keywords: dialogue systems, entity recognition, knowledge graphs

1. Introduction

When building dialogue systems, reliable named entity recognition (NER) and named entity linking (NEL) are vital to understanding user intent, especially if these dialogue systems are open domain and intended to support conversations on any topic. In designing our open domain social bot, Slugbot (Bowden et al., 2017b), for the 2017 Amazon Alexa Prize Contest (Ram et al., 2017), we discovered several challenges specifically associated with both NER and NEL. This paper discusses these challenges, and shows how we address them with Slugbot's Named Entity Recognition for Dialogue Systems (SlugNERDS), a tool designed for NER and NEL in open domain dialogue. Additionally we present two corpora, SlugEntityDB and SchemaActuator, which are based on over 10,000 real user conversations with the system. We perform an extensive analysis of our system and the corpora to identify important areas of future work.

NER and NEL have been actively researched topics for decades (Finkel and Manning, 2009; Ratinov and Roth, 2009; Ritter et al., 2011; Derczynski et al., 2015; Nitish Gupta and Roth, 2017). However, the resulting entity classification is often coarse and does not encode an ontology. For example, Stanford NER features only a small number of abstract entity types such as PERSON, LOCATION, ORGANIZATION, and MISC (Manning et al., 2014; Finkel and Manning, 2009); these categories don't provide enough information for dialogue interpretation and generalization. Although other resources such as that from Ratinov and Roth (2009) utilize additional external knowledge by extracting 30 gazetteers from both the web and Wikipedia, the entity types are still not as varied as we need, and the framework lacks a clear ontology. Furthermore, the alignment of classes between systems can be inconsistent as there is no universally shared taxonomy between them and the various data streams necessary to support open domain conversation (Bowden et al., 2017a). While Ling and Weld

(2012) attempt to address this by using 112 fine-grained entity types consistent with Freebase (Bollacker et al., 2008), Freebase is no longer maintained and recent inspection has shown it to be significantly incomplete (Dong et al., 2014). While the accuracy of these state of the art NER systems can be quite high, ranging between 80-90% on long text, on short informal text, such as tweets, accuracies drop to between 30-50% (Derczynski et al., 2015). Tweets are much more representative of the data we see from users in an interaction with a social bot than newswire data. Specifically, utterances tend to be short, and due to the open domain setting, relevant context is not guaranteed.

However, unlike tweets, a dialogue system must maintain a discourse model which can sustain multi-turn dialogue for the duration of the user interaction. In fact, it is clear that a dialogue system, which necessarily must operate in real time, has inherent challenges not present in other settings. Real time systems must be optimized such that they function without significant response delays between turns; such delays can be introduced by approaches that rely on running a machine learning classifier in real time. Additionally, new named entities are very commonly discussed in a social setting and need to be recognized as trends and current popular topics change: this requires systems to be constantly (nightly) retrained on newly annotated data. Moreover, there must be significant noise tolerance - in a social setting it is often the case that users will speak more informally; the spoken domain can be even more challenging as an entity resolution tool must also account for automatic speech recognition (ASR) misinterpretations.

2. Overview of Tools and Corpora

To address these challenges we present Slugbot's Named Entity Recognition for dialogue Systems (SlugNERDS), an NER and NEL tool which leverages the Google Knowledge Graph API in conjunction with the Schema.org taxonomy to identify known entities. The tool is optimized with respect to noisy open domain conversation and is able to per-

1	User: Let's talk about comic books. Who is your favorite character?	In open domain dialogue, the user can engage us on any topic.
2	System: I like [Moon Knight] _{Thing} because [he] _{Thing} can see the future sometimes.	Fictional entities, such as obscure comic book characters, are classified as <i>Thing</i> . Using our web based NEL, we could also summarize he is a <i>Fictional Super Hero</i> .
3	User: I guess that is a cool power, when did [he] _{Thing} first appear?	SlugNERDS performs anaphoric resolution to perform discourse level entity linking, in this case associating <i>he</i> with <i>Moon Knight</i> .
4	System: I think [Marvel Comics] _{Corporation} debuted [Moon Knight] _{Thing} in [werewolf by night] _{Thing} number 32	In SlugNERDS, correctly identifying <i>Moon Knight</i> allows us to perform web level entity linking - we have a link to a relevant Wikipedia article where we find more information about the character, specifically the series he first appeared in, <i>Werewolf by Night</i> .
5	User: you know I kind of like scary space movies . I saw [alien covenant] _{Movie} the other day	The <i>movies</i> keyword matches directly to an entity type, thus increasing the systems probability of picking the correct movie entity.
6	System: I thought [Alien Covenant] _{Movie} was scary and action packed.	Often times users are most interested in discussing trending entities, <i>Alien Covenant</i> had been released in the theaters recently.
7	User: That's cool, have you seen the new [Star Wars] _{MovieSeries} movie ?	Both the <i>movie</i> keyword and the <i>seen</i> action are useful to detect all <i>Movie</i> type entities, such as <i>MovieSeries</i> , <i>MovieClip</i> , and <i>Movie</i> .
8	System: I haven't had a chance yet. Which is your favorite in the series?	Differentiating between <i>Movie</i> and <i>MovieSeries</i> allows us to curate better follow up question.
9	User: I think my favorite [Star Wars] _{MovieSeries} movie is [revenge of the sith] _{Movie}	We just need the partial title, <i>revenge of the sith</i> , to correctly identify the entity <i>Star Wars: Episode III - Revenge of the Sith</i> .
10	System: Who is your favorite actor in [it] _{Movie} ?	We know that <i>Movies</i> have <i>actors</i> .
11	User: I'm watching the [Lakers] _{SportsTeam} play right now, which is funny because I hate [Los Angeles] _{City} , even though I live in [California] _{State} - best state in the [United States] _{Country} .	We note the various levels of granularity, instead of just a location tag we can detect each destination with more precision.
12	System: Whose your favorite player on the [Lakers] _{SportsTeam} ?	Using the Schema.org attributes, we know that <i>SportsTeams</i> have <i>athletes</i> .
13	User: [albert einstein] _{Person} is very smart [he] _{Person} is my favorite scientist. I wonder if [he] _{Person} would have liked reading [the hunger games] _{BookSeries} [mocking-jay] _{Book} .	Note that users can abruptly shift to a different context at any time. In this example, the <i>mockingjay</i> movies scored higher than the books initially, however once we take into account the expectation of <i>Book</i> type entities via the <i>reading</i> action, the correct entities are selected.

Table 1: Sample dialogue with entities highlighted in addition to SchemaActuators annotations.

form both discourse and web based entity linking. Table 1 represents an annotated conversation based on real interactions with Slugbot. While one could enhance the result further by utilizing the state of the dialogue system to set system expectations, we are interested in evaluating our tool without making any assumptions of the system using it.

To supplement SlugNERDS we also present SlugEntityDB, an annotated corpus which can be used to evaluate our system. This represents to our knowledge the first Schema.org entity type annotated corpus for this task. The SlugEntityDB contains 2100 samples, 500 taken directly from real user data collected by SlugBot (Bowden et al., 2017b) during the inaugural Amazon Alexa Prize contest, 1600 synthesized such that we can easily verify the richness of the

corpus. Since we are operating in the open domain, it is also very likely that adding synthesized data will result in entities which have never been seen by the system previously. In this dataset the utterances are formatted similar to the input which a spoken dialogue system would receive from a state of the art ASR system. Each utterance is annotated in tuples which includes the direct strong overlap, entity types, and full entity name as per the Google Knowledge Graph. Table 2 includes a sample of this dataset.

Additionally, we provide the SchemaActuators corpus, a partially hand annotated probabilistic mapping between actions/specific keywords/phrases and entity classifications (such as indicating *watch* or *seen* are related to *Movie* entities). An example of these mappings can be seen in Table 1,

Utterance	Annotation
the lord of the rings was my childhood	(lord of the rings, MovieSeries Thing, The Lord of the Rings)
my favorite star wars movie is probably revenge of the sith	(revenge of the sith, Movie Thing, Star Wars: Episode III Revenge of the Sith); (star wars, MovieSeries Thing, Star Wars)
i want to visit black mountain	(black mountain, TouristAttraction Mountain Place Thing, Black Mountain)

Table 2: Samples from the SlugEntityDB.

where appropriate entries in the corpus are colored to match the associated entity. Currently the verbs used in this corpus have been hand annotated. These seed verbs are then expanded using synonym relations from Wordnet (Miller, 1995) and other lexical resources. Using these verbs in addition to prepositional phrases we automatically generate a list of short phrases associated with specific entities, such as *arrive at* for the *City* entity. We use a similar process of automatic expansion on the entity type to generate a list of candidate keywords which can potentially indicate an entity, such as associating the *flick* and *film* keyword with the *Movie* entity. It is our belief that this corpus will lead to improved results as it allows us to better adjust our system’s expectations.

3. Tools and Methodology

In this section, we describe the tools and methodology we use to build SlugNERDS.

3.1. Google Knowledge Graph

The SlugNERDS tool primarily utilizes the Google Knowledge Graph API¹ to identify known entities. A query to the Google Knowledge Graph API returns a list of the N most probable entities which are sorted based on an associated base score. This score is provided by the API and is assumed to be a combination of contextual overlap and entity popularity. Each entity has substantial meta data including the full name of the entity, Schema.org type classification, and both a brief and long description with an associated Wikipedia article. Having automatic access to an associated Wikipedia article allows us to reliably solve the task of web based NEL contingent on successfully performing NER. An example query result can be seen in Figure 1.

Utilizing this API to identify known entities is ideal as it isn’t domain specific, and requires no training from users of the tool. While directly querying the API in real time may sound expensive, this tool has been successfully deployed in our real time conversational agent, Slugbot, without significantly inhibiting the user experience.

¹<https://www.google.com/intl/bn/insidesearch/features/search/knowledge.html>

```
{
  '@context': {
    ...
  },
  '@type': 'ItemList',
  'itemListElement': [
    {
      '@type': 'EntitySearchResult',
      'resultScore': 801.711243,
      'result': {
        'description': '2005 film',
        'url': ...,
        '@id': 'kg:/m/0fdv3',
        '@type': [
          'Movie',
          'Thing'
        ],
        'detailedDescription': {
          'articleBody': ...,
          'license': ...,
          'url': ...
        },
        'name': 'Star Wars: Episode III – Revenge of the Sith'
      }
    },
    ...
  ]
}
```

Figure 1: Google Knowledge Graph Search API Result For the Query *revenge of the sith*

3.2. Schema.org

The Google Knowledge Graph API classifies entities using the Schema.org² entity ontology. Schema.org is an effort to create a richer web infrastructure by proposing common MicroData for entities within a website. Millions of websites which contain rich structured data across an array of subjects - such as IMDb³, BestBuy⁴, BarnesNNobles⁵, and Yelp⁶ - have already adopted this MicroData in some form.

This not only enforces consistency across a multitude of data streams, but allows us to connect common entities to their related attribute types (such as *SportsTeam* → *athlete* → *Person* → *birthDate*), allowing the system to retrieve a large set of possible next topics, related facts, and associated entities. We can further expand on potential topics by utilizing the schema ontology to access properties of the entities higher up in the hierarchy.

3.3. SlugNERDS Pipeline

Figure 2 represents the general SlugNERDS pipeline. Our *Name Entity Recognition* consists of two standard modules, *Entity Segmentation* and *Entity Classification* (Ritter et al., 2011; Collins and Singer, 1999; Downey et al., 2007). Subsequently we perform *Entity Linking* on the recognized entity. We will examine this process with the following example utterance: *"I think my favorite star wars is revenge of the sith"*. Please note that punctuation is not included in the user utterance.

3.3.1. Entity Segmentation

In order to refine our list of candidate strings to query we must break our text into reasonable chunks. We utilize a

²<https://www.schema.org>

³<https://www.imdb.com>

⁴<https://www.bestbuy.com>

⁵<https://www.barnesandnoble.com>

⁶<https://www.yelp.com>

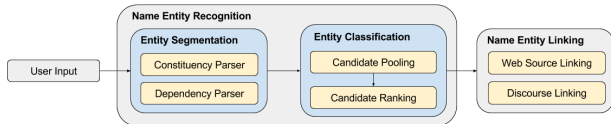


Figure 2: The SlugNERDS pipeline.

two pass approach to maximize our recall with a decent number of candidates. We maximize recall because a properly structured dialogue system will be able to pick the contextually relevant entity for follow-up questions and ignore extraneous entities which may have been misclassified. Furthermore, through empirical evaluation we have concluded that Slugbot missing an entity can be more detrimental to a conversation than over-classifying entities. First, we construct a constituency tree using Stanford CoreNLP (Manning et al., 2014) and build our candidate pool by collapsing each of the noun phrases, verb phrases, and sentences in the tree. A sample constituency tree can be seen in Figure 3. Additionally, we collapse sequential noun clusters from the dependency parse which have not yet been associated with an entity to create a secondary pool of candidates, so as to include more candidate strings that are ignored by shallow parsing (Ritter et al., 2011). We exclude single pronouns such as *I* and *me* unless they seem extremely contextually relevant, such as in the case of *It*⁷, when discussing movies with the user.

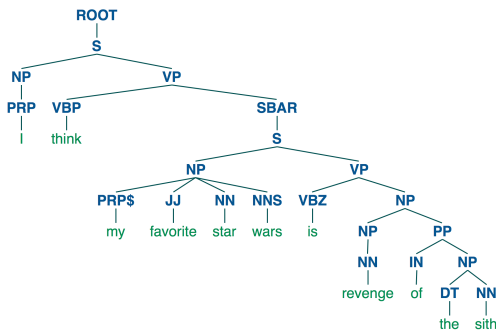


Figure 3: Constituency Tree

3.3.2. Entity Classification

Once we collect the candidate phrases in the Segmentation phase, we then query each of these phrases using the Google Knowledge Graph API and collect the top 6 relevant entities in *Candidate Pooling*. Here 6 is an empirically derived value which represents a good range of possible candidates without including excessive candidates. Sometimes it is possible for the entity candidate returned by the query to have the same exact title with different entity types. For example, the phrase *revenge of the sith* returns 5 candidates with the title *Star Wars: Episode III - Revenge Of The Sith* each with a different entity type (*Movie*, *Video Game*, *Book*, *MusicAlbum*, and *BookSeries*). The base score can vary significantly between these versions, for example the base score for the *Movie* entity is 795.59, while it is 138.05

⁷*It* is currently a very popular horror movie.

for the *BookSeries* entity - this is assumed to be due to the popularity of the movie vs. the book series.

<p>User Utterance: I think my favorite <i>Star Wars</i> movie is <i>revenge of the sith</i>.</p> <p>Candidate: Star Wars: Episode III - Revenge of the Sith [<i>'Movie', 'Thing'</i>]</p> <p>Candidate Initial Vector: [1, 1, 1, 1, 1, 1, 1, 1]</p> <p>Candidate Overlap: [1, 1, 0, 0, 0.11, .75, .75, 1]</p>

Table 3: Example word vector translation.

Since the entities returned by the Google Knowledge Graph may not be an exact match to our query, it gives us more flexibility, while introducing some noise. Furthermore, if we are expecting a user to talk about certain entity types according to the context, as discussed when introducing the SchemaActuators corpus in Section 2., we place increased value on certain entities while penalizing others. Thus, we perform our scoring algorithm to maximize our performance in the *Candidate Ranking* phase.

First, we create a word overlap vector indicating the total overlap. Because we expect noise in the user's utterance, we allow for unexpected words to be inserted in the middle of the phrase realizing the entity with a distance based penalty. Table 3 gives an example of the word vector. The word vector is initialized as 1 for each candidate token, for instance, *Star Wars: Episode III - Revenge of the Sith* has 8 tokens, thus the initial vector is [1,1,1,1,1,1,1,1]. Then we compare the candidate tokens with the user utterance. The element of the vector will be penalized if tokens are missing or extra tokens are detected. For example, the element 3 and element 4 are set to zero because there is two missing words (*Episode III*), resulting in the 5th element of the vector receiving a penalty of $\frac{1}{3}$, which is the distance between the matching tokens. There are also two extra words (*movie is*) in the user utterance, thus an additional penalty of $\frac{1}{3}$ is applied to the 5th element, which yields a 0.11 score. We also note the imposed penalty on stop words by multiplying the element with 0.75 - this helps to increase our precision as we try to prevent correctly positioned stop words from accidentally forming an entity. The 6th and 7th elements (*of the*) are the stop word examples. Equation 1 demonstrates how we use this word vector to alter the score of an entity.

$$overlap_score = \frac{sum(word_vector)}{(\frac{1}{word_vector.count(0) + 1}) * base_score} \quad (1)$$

Secondly, we account for any entities we are expecting based on early inspection of the utterance's context:

$$can_score = overlap_score * (num_matches + 1) \quad (2)$$

Once the entities are all scored, we rerank our list and consider only the top ranked entity for each node, while also pruning away nodes whose top scoring entity was less than a certain threshold (empirically driven). Finally, we merge overlapping nodes who have candidates. For example, if both *revenge of the sith* and a child node *revenge* have the

same entity as their top scoring entity, we will merge these two nodes or remove lower ranking conflicts.

In our last stage, we sync the query/candidate to our internal discourse state representation. In our example, two entities are extracted, mapping *star wars* to entity type *MovieSeries*, and *revenge of the sith* to *Star Wars: Episode III - Revenge Of The Sith* with entity type *Movie*.

3.3.3. Entity Linking

Named Entity Linking is primarily encapsulated in two phases, *Web Source Linking* and *Discourse Linking*. With *Web Source Linking* we are interested in linking a known entity to existing resources on the web while *discourse linking* is focused on linking each mention of the entity within the input to the same discourse entity in our internal representation (Brennan et al., 1987; Walker et al., 1997).

As mentioned in Section 3.1., the Google Knowledge Graph query returns a Wikipedia article associated with the entity. We can further increase our web based linking by utilizing the fact that a large number of popular websites use the Schema.org MicroData, allowing us to easily target relevant sources for information extraction. Finally, through empirical examination, we note that pairing the entity type with the precise entity name as provided in the query will allow for easy subsequent queries to large databases such as YAGO (Rebele et al., 2016) or DBpedia (Auer et al., 2007).

Finally, our tool uses an augmented version of the Stanford Coreference Annotator (Manning et al., 2014) to perform *Discourse Linking*.

4. Evaluation

4.1. NER Results

Our SlugNERDS tool was originally developed and utilized in the 2017 Alexa Prize Competition for SlugBot, which scored in the top 25% of competing social bots. For detailed system evaluation, we present a set of experiments to evaluate our SlugNERDS tool and SchemaActuator corpus, independently from SlugBot. We evaluate SlugNERDS using (1) the base scores using only text segmentation with the Google Knowledge Graph, (2) augmenting the scores using our scoring algorithm, (3) using the SchemaActuator actions to increase context, (4) using the SchemaActuator keywords to increase context, (5) an ensemble approach. Moreover we test these configurations with two different scoring thresholds (the minimum score required for an entity be accepted as the correct class). The two thresholds we will test are (1) 150, an empirically driven value deemed to be a reasonable threshold during the development of SlugBot and (2) 0, no threshold at all.

Table 4 presents the results of our NER experiments using the SchemaEntityDB. Since Stanford NER is still commonly used in many state of the art open domain conversational systems, we use it as our baseline system. Our other experiments include using just our text segment method, then adding SlugNERDS ranking, iteratively adding in verbs and nouns from our SchemaActuators corpus, and finally showing our ensemble method results (which merges all resources). Since Stanford provides coarse grain entity type such as PERSON, LOCATION, ORGANIZATION,

and MISC, we map the specific gold standard label to these four types for the evaluation.

We use accuracy, macro-f1, micro-f1, and weighted-f1, to have a better understanding of our system's performance. High accuracy entity detection will ensure a more satisfactory conversation, and prevent us from missing the topics of the conversation. The macro-f1 treats all the classes evenly, the micro-f1 accounts for label imbalance, and the weighted-f1 is a weighted macro-f1. Since detecting infrequent entities is also important for us, we evaluate the macro-f1, the micro-f1 and the weighted-f1 for different interests. Table 4 shows that Stanford yields a macro-F1 score as 0.079. As predicted, we see from these results that Stanford NER is not a suitable system to use when detecting entities in open domain discourse. Our text segmentation model has an accuracy of 0.751, macro-F1 of 0.785, and weighted-F1 of 0.831. Our SlugNERDS models has a better accuracy 0.777 and a slightly worse macro-F1 of 0.746, but a better weighted-F1 0.85. After utilizing the SchemaActuator corpus, the macro-F1 are improved slightly, though the weighted-F1 is almost the same, which might due to the reason that the extra Act.Verbs and Act.Nouns are able to detect the edge cases which are infrequent entities. We believe that when integrated with a real dialogue system, the increase in contextual knowledge from our dialogue manager will yield further increase performance.

In our original experiment we see that a low threshold doesn't need contextual information to classify entities. However after examining the results it's clear that while more correct entities are being classified, there is also increased levels of over-classification which can be detrimental to the system. Therefore we aim to increase our accuracy with an with an empirically derived threshold of 150. In this experiment our model continued to outperformed the stanford baseline and the text segmentation results were the worst configuration for SlugNERDS. We can also see by comparing Table 4 and Table 5 that as we increase this threshold, we increase the impact of adding contextual information to SlugNERDS. This implies the importance of encoding contextually relevant data in our model while also increasing our tolerance of noise. In future work we aim at further analyzing our distribution of annotated entities to pick an optimized threshold value.

4.2. Distribution of Detected Entities

Figure 4 shows the distribution of Stanford Entities that are detected. The NUMBER, DATE, TIME, SET, ORDINAL, DURATION, MONEY, and PERCENT are not very useful to our system as we are interested primarily in known entities. Therefore we are most interested in the PERSON, ORGANIZATION, LOCATION, and MISC types - however, all of these entity types detected much less frequently than they appear in our annotated data, as can be seen in Figure 6. Figure 5 shows the top 15 most frequent entities that caught our SlugNERDS Models. We can see that our SlugNERDS model is able to successfully detect large amounts of various conversational entities such as Movie, Book, and MusicRecording, while also maintaining a very similar distribution as seen in our annotated data.

	Accuracy	macro-F1	micro-F1	weighted-F1
Stanford Baseline	.029	.079	.029	.052
Text Segmentation Only	.751	.785	.722	.831
SlugNERDS Scoring	.777	.746	.749	.850
SlugNERDS Scoring + Act_Verbs	.770	.752	.747	.849
SlugNERDS Scoring + Act_Nouns	.773	.762	.745	.848
Ensemble	.763	.747	.742	.845

Table 4: NER Results for the Threshold 0 experiment.

	Accuracy	macro-F1	micro-F1	weighted-F1
Stanford Baseline	.029	.079	.029	.052
Text Segmentation Only	.521	.656	.511	.663
SlugNERDS Scoring	.553	.624	.546	.692
SlugNERDS Scoring + Act_Verbs	.604	.655	.598	.735
SlugNERDS Scoring + Act_Nouns	.579	.619	.571	.714
Ensemble	.592	.634	.586	.726

Table 5: NER Results for the Threshold 150 experiment.

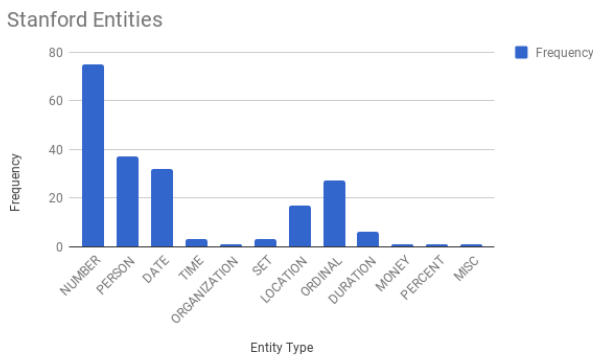


Figure 4: Distribution entity types detected by the Stanford NER.

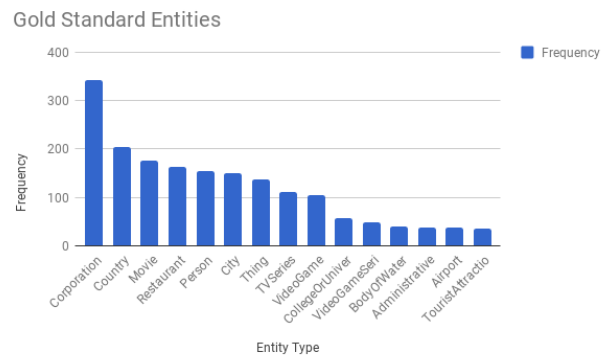


Figure 6: Distribution of top 15 entity types in the annotated data.

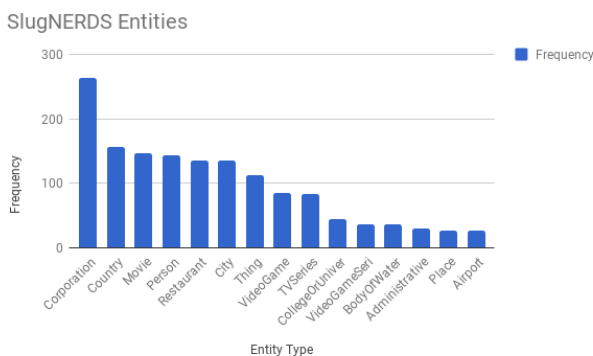


Figure 5: Distribution of top 15 entity types that our model detected.

4.3. Error Analysis

Here we will discuss different errors we noticed while analyzing the results of the experiments described above. Primarily, we see three different classes of errors; insufficient contextual information, natural language understanding de-

ficiencies, and the difficulty of encapsulating every possible entity within the Google Knowledge Graph.

Table 6 demonstrates several errors which are the result of insufficient contextual information. While this evaluation was meant to analyze SlugNERDS as a stand-alone tool, it is clear that these errors can be resolved by enhancing the contextual knowledge we pass to our tool. For example, while there are no lexemes which indicate *lord of the rings* in Sample 6a is referring to the MovieSeries, the state of our dialogue system could inform us that we are discussing movies, rather than books. As in Sample 6b, we see that it common in colloquial speech to refer to an element of a series by it's common root word. For example *halo* in the expression *let's play halo* is valid for *halo combat evolved*, *halo 2*, *halo wars*, and various other titles in the series. Disambiguating this is not a trivial problem and requires a significantly more rich state than the previous example. We will note here that it is possible to also leverage the granularity of our Schema.org entity types to recognize we are talking about a VideoGameSeries and clarify the specific VideoGame with the user. Finally, Sample 6c

demonstrates how a lack of user meta-data, in this case location and frequency, results in incorrect classifications. A dialogue system can represent this contextual data by referencing a user model which may indicate that since the user lives in California, they are likely referring to San Jose, California, rather than San Jose, Costa Rica.

Table 7 demonstrates two areas in which adding an additional layer of NLU would yield increased performance. Specifically we notice in Sample 7b that abbreviations cause difficulty in classifying entities. While *ucsd* will in fact return the correct entity as a potentially candidate from the Google Knowledge Graph, there is no lexical overlap, resulting in a false classification. By adding an additional layer of NLU which is able to expand abbreviated entities we would see an increase in performance. As described previously, the *SlugEntityDB* annotated corpora was designed assuming the input provided is from a spoken dialogue system - meaning it will suffer from the limitations of state of the art automatic speech recognition. More specifically, our utterances have no punctuation, capitalization, or non-alpha-numeric symbols. While we have already resolved this partially such that utterances like *x man* will correctly map to *X-man*, the lack of punctuation and capitalization can lead to inaccurate results from our parser - this can directly alter the queries which are sent to the Google Knowledge Graph. This can be seen in Sample 7a, where the best possible query for the original utterance was *love gordon ramsey* which does not yield the correct entity. After manually capitalizing the entity name, our best possibly query became *Gordon Ramsey*, which resulted in the correct entity.

Finally, as there exists an infinite amount of entities, it is reasonable to surmise that the Google Knowledge Graph, while quite robust, is incomplete. For example, querying either *xbox 360* or *windows* both return the *Microsoft* entity has the best possible candidate. While not necessarily common, this can be a difficult problem to deal with, especially when an entire category of entities, such as operating systems or video game consoles, are missing from the Knowledge Graph.

5. Conclusion and Future Work

In this paper we have presented *SlugNERDS*, an NER and NEL tool which is optimized with the respect to the challenges that are specific to open domain conversation. We have also presented two relevant corpora, one being the first dataset of its kind to be annotated with Schema.org named entity types in addition to the SchemaActuator corpus - a mapping of actions and keywords to their respective Schema.org entity types.

To our knowledge our system is the only one to utilize the Schema.org entity types for entity classification, but we plan to compare our system more extensively to other existing NER systems in future work, which may allow us to improve our system. One is T-NER System (Ritter et al., 2011)⁸, which is optimized for NER in Tweets, as discussed in Section 1.. We are also interested in comparing against a state of the art neural NEL model (Nitish Gupta

Utterance	Correct Entity	Predicted Entity
the lord of the rings was my childhood	MovieSeries	BookSeries

(a) Without sufficient contextual knowledge, it is difficult to differentiate between entities with the exact same title.

Utterance	Correct Entity	Predicted Entity
halo has been dead for a while now	VideoGame, Halo: Combat Evolved	VideoGameSeries, Halo)

(b) Entities which are nested within a series can also be hard to detect, such as here when the user is talking about the first Halo game vs. the Halo series.

Utterance	Correct Entity	Predicted Entity
sacramento airport was pretty busy	Sacramento International Airport	Sacramento Airport
i did not know san jose is a capital	San Jose, California	San Jose

(c) Without contextual information about the user, it's difficult to differentiate between San Jose refers to San Jose California vs. San Jose Costa Rica. Similarly, "common sense" indicates that "Sacramento Airport" refers to Sacramento International Airport, rather than the much smaller Sacramento Airport.

Table 6: Common errors stemming from insufficient contextual knowledge.

Utterance	Correct Entity	Predicted Entity
i love gordon ramsay	Person	None
i love Gordon Ramsay	Person	Person

(a) A lack of capitalization and punctuation leads to parser errors, resulting in missed queries.

Utterance	Correct Entity	Predicted Entity
ucsd is number 23	University of California, San Diego	None

(b) Mismatch between the string in the utterance vs the actual entity name, as in the case with an abbreviation, causes classification to fail.

Table 7: Common errors stemming from gaps in our natural language understanding pipeline.

and Roth, 2017)⁹. This model more proactively tackles the open domain problem, and can include updated knowledge of existing entities without retraining. There are clearly many areas in which we can further evaluate and improve the performance of *SlugNERDS*. Specifically, we plan to evaluate the performance of *SlugNERDS* with an emphasis on the impact of contextual information provided as a result of our dialogue system's state tracking capabilities. We also plan to improve our natural language understanding pipeline by investigating methodologies for improving in-

⁸https://github.com/aritter/twitter_nlp

⁹<https://nitishgupta.github.io/neural-el/>

put text quality, such as automatic capitalization and punctuation insertion. Finally, we are currently investigating how reinforcement learning can be applied to our existed conversational data to increase our models contextual insight.

6. References

- Auer, S., Bizer, C., Kobilarov, G., Lehmann, J., Cyganiak, R., and Ives, Z. (2007). Dbpedia: A nucleus for a web of open data. In *Proceedings of the 6th International The Semantic Web and 2Nd Asian Conference on Asian Semantic Web Conference, ISWC'07/ASWC'07*, pages 722–735, Berlin, Heidelberg. Springer-Verlag.
- Bollacker, K., Evans, C., Paritosh, P., Sturge, T., and Taylor, J. (2008). Freebase: A collaboratively created graph database for structuring human knowledge. In *Proceedings of the 2008 ACM SIGMOD International Conference on Management of Data, SIGMOD '08*, pages 1247–1250, New York, NY, USA. ACM.
- Bowden, K. K., Oraby, S., Wu, J., Misra, A., and Walker, M. (2017a). Combining search with structured data to create a more engaging user experience in open domain dialogue. *SCAI 2017*.
- Bowden, K. K., Wu, J., Oraby, S., Misra, A., and Walker, M. (2017b). Slugbot: An application of a novel and scalable open domain socialbot framework. In *Alexa Prize Proceedings*.
- Brennan, S. E., Friedman, M. W., and Pollard, C. J. (1987). A centering approach to pronouns. In *Proc. 25th Annual Meeting of the ACL, Stanford*, pages 155–162.
- Collins, M. and Singer, Y. (1999). Unsupervised models for named entity classification. In *In Proceedings of the Joint SIGDAT Conference on Empirical Methods in Natural Language Processing and Very Large Corpora*, pages 100–110.
- Derczynski, L., Maynard, D., Rizzo, G., van Erp, M., Gorrell, G., Troncy, R., Petrak, J., and Bontcheva, K. (2015). Analysis of named entity recognition and linking for tweets. *Information Processing and Management*, 51(2):32 – 49.
- Dong, X., Gabrilovich, E., Heitz, G., Horn, W., Lao, N., Murphy, K., Strohmman, T., Sun, S., and Zhang, W. (2014). Knowledge vault: A web-scale approach to probabilistic knowledge fusion. In *Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '14*, pages 601–610, New York, NY, USA. ACM.
- Downey, D., Broadhead, M., and Etzioni, O. (2007). Locating complex named entities in web text. In *IJCAI 2007, Proceedings of the 20th International Joint Conference on Artificial Intelligence, Hyderabad, India, January 6-12, 2007*, pages 2733–2739.
- Finkel, J. R. and Manning, C. D. (2009). Joint parsing and named entity recognition. In *Proceedings of Human Language Technologies: The 2009 Annual Conference of the North American Chapter of the Association for Computational Linguistics, NAACL '09*, pages 326–334, Stroudsburg, PA, USA. Association for Computational Linguistics.
- Ling, X. and Weld, D. S. (2012). Fine-grained entity recognition. In *Proceedings of the Twenty-Sixth AAAI Conference on Artificial Intelligence, AAAI'12*, pages 94–100. AAAI Press.
- Manning, C. D., Surdeanu, M., Bauer, J., Finkel, J., Bethard, S. J., and McClosky, D. (2014). The Stanford CoreNLP natural language processing toolkit. In *Association for Computational Linguistics (ACL) System Demonstrations*, pages 55–60.
- Miller, G. A. (1995). Wordnet: A lexical database for english. *Commun. ACM*, 38(11):39–41, November.
- Nitish Gupta, S. S. and Roth, D. (2017). Entity linking via joint encoding of types, descriptions, and context. In *EMNLP*.
- Ram, A., Prasad, R., Khatri, C., Venkatesh, A., Gabriel, R., Liu, Q., Nunn, J., Hedayatnia, B., Cheng, M., Nagar, A., King, E., Bland, K., Wartick, A., Pan, Y., Song, H., Jayadevan, S., Hwang, G., and Pettigru, A. (2017). Conversational ai: The science behind the alexa prize. In *Alexa Prize Proceedings*.
- Ratinov, L. and Roth, D. (2009). Design challenges and misconceptions in named entity recognition. In *Proceedings of the Thirteenth Conference on Computational Natural Language Learning, CoNLL '09*, pages 147–155, Stroudsburg, PA, USA. Association for Computational Linguistics.
- Rebele, T., Suchanek, F., Hoffart, J., Biega, J., Kuzey, E., and Weikum, G., (2016). *YAGO: A Multilingual Knowledge Base from Wikipedia, Wordnet, and Geonames*, pages 177–185. Springer International Publishing, Cham.
- Ritter, A., Clark, S., Mausam, and Etzioni, O. (2011). Named entity recognition in tweets: An experimental study. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing, EMNLP '11*, pages 1524–1534, Stroudsburg, PA, USA. Association for Computational Linguistics.
- Marilyn A. Walker, et al., editors. (1997). *Centering in Discourse*. Oxford University Press.