

# Transformation-based Learning for Semantic parsing

*F. Jurčiček, M. Gašić, S. Keizer, F. Mairesse, B. Thomson, K. Yu, and S. Young*

Engineering Department, Cambridge University, Trumpington Street, Cambridge, CB2 1PZ, UK

{fj228, mg436, sk561, farm2, brmt2, ky219, sjy}@eng.cam.ac.uk

## Abstract

In this paper, we present a semantic parser that transforms an initial semantic hypothesis into correct semantics by using an ordered list of rules. These rules are learnt automatically from a training corpus with no prior linguistic knowledge and no alignment between words and semantic concepts. We show that this parser is competitive with respect to the state-of-the-art semantic parsers on the ATIS and TownInfo tasks.

**Index Terms:** spoken language understanding, semantics, natural language processing, transformation-based learning

## 1. Introduction

The goal of semantic parsing is to map natural language to a formal meaning representation - semantics. Such semantics can be either defined by a grammar, e.g. LR grammar for GeoQuery domain [1], or by frames and slots, e.g. TownInfo domain [2]. Table 1 shows an example of the frame and slot semantics from the ATIS dataset [3]. Each frame has a goal and set of slots. Each slot is composed of a slot name, e.g. “from.city”, and slot value, e.g. “Washington”. As dialogue managers commonly use semantics in the form of frames and slots [4, 5], our approach learns to map directly from natural language into frame and slot semantics.

A dialogue system needs a semantic parser which is accurate and robust, easy to build, and fast. First, this paper presents a parsing technique that performs comparably to state-of-the-art semantics parsers and it can handle ill formed utterances. Second, it does not need any handcrafted linguistic knowledge. Third, it can learn from data which is not semantically annotated at the word-level. Finally, it learns a compact set of rules that can be used for realtime semantic parsing.

In our approach, we adapt Transformation-Based Learning (TBL) [6] to the problem of semantic parsing. We attempt to find an ordered list of transformation rules which iteratively improve the initial semantic annotation. In each iteration, a transformation rule corrects some of remaining errors in the semantics. To handle long-range dependencies between words, we experiment with features extracted from dependency parse trees provided by the RASP syntactic parser [7].

In the next section, we describe previous work on mapping natural language into a formal meaning representation. Section 3 presents an example of TBL semantic parsing, discusses templates for transformation rules, and describes the learning process. Section 4 compares the TBL parser to the previously developed semantic parsers on the ATIS [3] and TownInfo [2] domains. Finally, Section 5 concludes this work.

what are the lowest airfare from Washington DC to Boston

GOAL	=	airfare
airfare.type	=	lowest
from.city	=	Washington
from.state	=	DC
to.city	=	Boston

Table 1: Example of frame and slot semantics from the ATIS [3] dataset.

## 2. Related work

In Section 4, we compare the performance of our method with four existing systems that were evaluated on the same dataset. First, the Hidden Vector State (HVS) technique has been used to model an approximation of a pushdown automaton with semantic concepts as non-terminal symbols [8, 9]. Second, Probabilistic parser using Combinatory Categorical Grammar (PCCG) has been used to map utterances to lambda-calculus [10]. This technique produces state-of-the-art performance on the ATIS dataset. However, apart from using the lexical categories (city names, airport names, etc) readily available from the ATIS corpus, this method also needs a considerable number of hand-crafted entries in their initial lexicon. Third, Markov Logic Networks (MLN) have been used to extract slot values by combining probabilistic graphical models and first-order logic [11]. In this approach, weights are attached to first-order clauses which represent the relationship between slot names and their values. Such weighted clauses are used as templates for features of Markov networks. Finally, in Semantic Tuple Classifier (STC) support vector machines have been used to build semantic trees by recursively calling classifiers that predict fragments of the semantic representation from n-gram features [2].

Also, there is a large amount of research that is not directly comparable because of either different corpora or different meaning representation. Transformation techniques have been used to sequentially rewrite an utterance into semantics [1]. However, our approach differs in the way the semantics is constructed. Instead of rewriting an utterance, we transform an initial semantic hypothesis. As a result, the words in the utterance can be used several times to trigger transformations of the semantics. Support vector machines and tree kernels [12] have been used to integrate knowledge contained in the manually annotated dependency trees to capture long-range relationships between words.

## 3. Transformation-based parsing

This section describes the transformation-based parser. First, we give an example how the parsing algorithm works. Second, we detail locality constraints on the transformation rules. Third,

we describe the templates used to generate rules for the learning process. Fourth, we describe features capturing long-range dependencies. Finally, the learning process is detailed.

### 3.1. Example of Parsing

The semantic parser transforms an initial semantic hypothesis into correct semantics by applying transformations from a list of rules. Each rule is composed of a trigger and a transformation and a trigger initiates transformation of a hypothesis. This section demonstrates the parsing process on the example: “*find all the flights between Toronto and San Diego that arrive on Saturday*”

First, the goal “flight” is used as the initial semantics because it is the most common goal in the ATIS dataset. As a result, the initial semantics is as follows:

GOAL = flight

Second, the rules which triggers match the utterance and the hypothesis are sequentially applied. Generally, the rules can add, delete, or substitute slots.

trigger	transformation
“between toronto”	add slot “from.city=Toronto”
“and san diego”	add slot “to.city=San Diego”
“saturday”	add slot “departure.day=Saturday”

After applying the transformations, we obtain the following semantic hypothesis:

GOAL = flight  
from.city = Toronto  
to.city = San Diego  
departure.day = Saturday

As the date and time values are associated with the “departure.\*” slots most of the time in the ATIS dataset, the parser learns to associate them with the “departure.\*” slots. The incorrect classification of the word “Saturday” is a result of such a generalisation. However, the TBL method learns to correct its errors. Therefore, the parser also applies the error correcting rules at a later stage. For example, the following rule corrects the slot name of the slot value “Saturday”.

trigger	transformation
“arrive”	substitute slot from “departure.day=*” to “arrival.day=*”

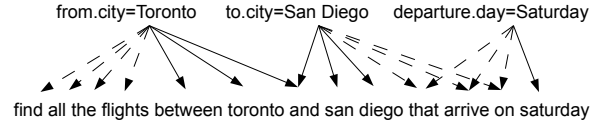
In this case, we substitute the slot name with the correct name, to produce the following semantic hypothesis:

GOAL = flight  
from.city = Toronto  
to.city = San Diego  
arrival.day = Saturday

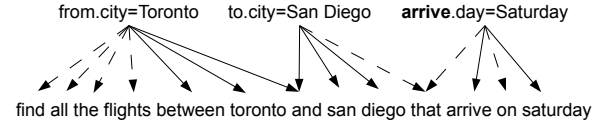
### 3.2. Locality constraints

So far no relationship between slots and their lexical realisation was considered. For example, before we perform the substitution of the slot “departure.day” to “arrival.day”, we should test whether the word “arrive” is near the slot’s lexical realisation. The reason is that we do not want to trigger the substitution of the slot “from.city=Toronto” to “to.city=Toronto”. This could happen because the parser has also learnt the following rule:

trigger	transformation
“arrive”	substitute slot from “from.city=*” to “to.city=*”



(a) alignment after the first set of rules



(b) alignment after application of the substitution rule

Figure 1: Alignment between the words and the slots in the example utterance.

One way to handle this problem is to constrain triggers performing substitution to be activated only if they are in the vicinity of the slot’s lexical realisation. We track the words from the utterance that were used in triggers. Every time we apply a transformation of a slot, we store links between the words which triggered the transformation and the target slot. Such links are referred as “direct alignment”.

In Figure 1 (a), we see the alignment between the words and the slots in the example utterance after applying the first set of rules. The full arrows denote direct alignment created by the add-transformations. Because no rules were triggered by the words “find all the flights” and “that arrive on”, those words could not be aligned directly to any of the slots. Therefore, we have to derive such alignment (see Figure 1 (a) dashed arrows). A word is aligned to a slot if the alignment does not cross any direct alignment. In Figure 1 (a), the phrase “find all the flights” can be aligned to the slot “from.city=Toronto” only (dashed arrows). The phrase “that arrive on” can be aligned to two slots “to.city=San Diego” and “departure.day=Saturday”.

In Figure 1 (b), we see the alignment after applying the substitution from Section 3.1. We can see a change in the alignment of the phrase “that arrive on”. First, the word “arrive” is aligned to the slot “arrival.day=Saturday”. Second, the word “on” must be aligned to the same slot as the word “arrive”. There is no change in the alignment of the word “that”.

### 3.3. Rule templates

In the previous sections, we presented several rules which add, substitute, or delete slots. These rules were selected by a learning algorithm from a large set of potential rules. Such rules are generated from a set of templates for triggers and transformations.

A trigger tries to match an input utterance, an output semantics, or both. In our method, a trigger contains one or more conditions as follows: the utterance contains n-gram  $N$ , the utterance contains skipping bigram  $B^1$ , the goal equals to  $G$ , and the semantics contains slot  $S$ . If a trigger contains more than one condition, then all conditions must be satisfied. In our method, we use unigrams, bigrams, trigrams and skipping bigrams, which can skip up to 3 words.

A transformation performs one of the following operations: substitute a goal to  $G$ , add a slot  $S$ , delete a slot  $S$ , and substitute a slot  $S$ . A substitution transformation can substitute a whole slot, a slot name, or a slot value.

<sup>1</sup> In a skipping bigram, one, two or three words are skipped between two words.

### 3.4. Improving the disambiguation of long-range dependencies

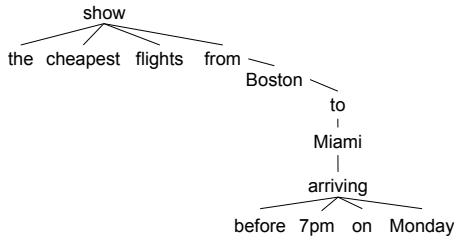


Figure 2: Dependency tree of the utterance "show the cheapest flights from Boston to Miami arriving before 7pm on Monday".

Besides simple n-grams and skipping bigrams, more complex lexical features can be used. Kate [12] used manually annotated dependency trees to capture long-range relationships between words. In a dependency tree, each word is viewed as the dependant of one other word, with the exception of the root. Dependency links represents gramatical relationships between words. Kate showed that word dependencies significantly improve semantic parsing because long-range dependencies from an utterance tend to be local in a dependency tree. For example, the words "arriving" and "Monday" are neighbours in the dependency tree but they are four words apart in the utterance (see Figure 2).

Instead of using manually annotated word dependencies [12], we used dependencies provided by the RASP dependency parser [7]. New n-gram features were generated in which a word history is given by links between words. For example, the algorithm would generate bi-gram ('arriving', 'Monday') for the word "Monday". We however note that RASP was used off-the-shelf and more accurate dependencies should be obtained by adapting the RASP parser to the ATIS and TownInfo domains.

### 3.5. Learning

The main idea behind transformation-based learning is to learn an ordered list of rules which incrementally improve an initial semantic hypothesis (see the algorithm in Figure 3)<sup>2</sup>. The initial assignment is made based on simple statistics - the most common goal is used as initial semantics. The learning is conducted in a greedy fashion, and at each step the algorithm chooses the transformation rule that reduces the largest number of errors in our hypothesis. Errors include goal substitutions, slot insertions, slot deletions, and slot substitutions. The learning process stops when no rule that improves the hypothesis beyond the pre-set threshold can be found. Note that no alignment between words and semantic concepts is needed.

As in the previous work [2, 8, 10, 11], we make use of a database with lexical realisations of some slots, e.g. city and airport names. Since the number of possible slot values for each slot is usually very high, the use of a database results in a more robust parser. In our method, we replace lexical realisations of slot values with category labels before parsing, e.g. "i want to fly from CITY". After parsing we use a deterministic algorithm to recover the original values for category labels, which is detailed in [2].

<sup>2</sup>The list of rules must be ordered because each learnt rule corrects errors remaining after application of preceding rules.

1. ASSIGN INITIAL SEMANTICS TO EACH UTTERANCE
2. REPEAT AS LONG AS THE NUMBER OF ERRORS ON THE TRAINING SET DECREASES
  - (A) GENERATE ALL RULES WHICH CORRECT AT LEAST ONE ERROR IN THE TRAINING SET
  - (B) MEASURE THE NUMBER OF CORRECTED ERRORS BY EACH RULE
  - (C) SELECT THE RULE WITH THE LARGEST NUMBER OF CORRECTED ERRORS
  - (D) APPLY THE SELECTED RULE TO THE CURRENT STATE OF THE TRAINING SET
  - (E) STOP IF THE NUMBER OF CORRECTED ERRORS IS SMALLER THAN THRESHOLD T.

Figure 3: Rule learning algorithm.

## 4. Evaluation

In this section, we evaluate our parser on two distinct corpora, and compare our results with state-of-the-art techniques and a handcrafted Phoenix parser [13].

### 4.1. Datasets

In order to compare our results with previous work [2, 8, 10, 11], we apply our method to the ATIS dataset [3]. We use 5012 utterances for training, and the DEC94 dataset as development data. As in previous work, we test our method on the 448 utterances of the NOV93 dataset, and the evaluation criterion is the F-measure of the number of reference slot/value pairs that appear in the output semantics (e.g., from.city = New York). He & Young detail the test data extraction process in [8].

Our second dataset consists of tourist information dialogues in a fictitious town (TownInfo). The dialogues were collected through user trials in which users searched for information about a specific venue by interacting with a dialogue system in a noisy background. The TownInfo training, development, and test sets respectively contain 8396, 986 and 1023 transcribed utterances. The data includes the transcription of the top hypothesis of a speech recogniser, which allows us to evaluate the robustness of our models to recognition errors (word error rate = 34.4%). We compare our model with the STC parser [2] and the handcrafted Phoenix parser [13]. The Phoenix parser implements a partial matching algorithm that was designed for robust spoken language understanding.

### 4.2. Results

The results for both datasets are shown in Table 2. The model accuracy is measured in terms of precision, recall, and F-measure (harmonic mean of precision and recall) of the slot/value pairs. Both slot and value must be correct to count as a correct classification.

Results on the ATIS dataset show that our method (F-measure = 95.74% ) is competitive with respect to the Zettlemoyer & Collins' PCCG model [10] (95.9%). Note that this PCCG model makes use of a considerably large number of handcrafted entries in their initial lexicon. In addition, our method outperforms the STC [2], HVS [8] and MLN [11] parsers. Concerning the TownInfo dataset, Table 2 shows that TBL produces 87.82% of F-measure, which represents a 3.28% improvement over the handcrafted Phoenix parser, while being competitive with the STC model - TBL's performance is only

Parser	Prec	Rec	F
<b>ATIS dataset with transcribed utterances:</b>			
TBL	96.37	95.12	95.74
PCCG	95.11	96.71	95.9
STC	96.73	92.37	94.50
HVS	-	-	90.3
MLN	-	-	92.99
<b>TownInfo dataset with transcribed utterances:</b>			
TBL	96.05	94.66	95.35
STC	97.39	94.05	95.69
Phoenix	96.33	94.22	95.26
<b>TownInfo dataset with ASR output:</b>			
TBL	92.72	83.42	87.82
STC	94.03	83.73	88.58
Phoenix	90.28	79.49	84.54

Table 2: Slot/value precision (Prec), recall (Rec) and F-measure (F) for the ATIS and TownInfo datasets.

Parser	Prec	Rec	F
<b>ATIS development dataset:</b>			
TBL	93.95	93.70	93.82
No locality constrains	93.38	92.64	93.01
No dependency tree features	92.78	92.04	92.41

Table 3: Comparison of different aspects of the TBL method on the ATIS development dataset.

0.76% lower.

Table 3 shows contrast between the full system and the system with no features extracted from dependency trees and the system with no locality constraints. Experiments were carried out on the ATIS development dataset. The results show that if the dependency tree features are removed or the locality constraints are not used, the performance of our method degrades.

The learning time of the TBL parser is acceptable and the parsing process is fast. First, the learning time is about 24 hours on an Intel Pentium 2.8GHz for each dataset. The TBL parser generate up to 1M of transformation rules in each iteration; however, only fraction of rules has to tested because the search space can be efficiently organized [6]. Second, the TBL method is able to decode an utterance in about 6ms. Because the number of learnt rules is considerably low, a number of operations needed to parse a utterance is limited. There are 17 unique dialogue acts and 66 unique slots in the ATIS dataset and the total number of learnt rules is 372. This results in 4.5 rules per semantic concept on average. In the TownInfo dataset, we have 14 dialogue acts and 14 slots and the total number of learnt rules is 195. The average number of rules per semantic concept is 6.9. The number of semantic concepts per utterance is 5 on average.

## 5. Conclusion

This paper presents a novel application of TBL for semantic parsing. Our method learns a sequence of rules which iteratively transforms initial semantics into correct semantics. Our method was applied to two very different domains and it was shown that its performance is competitive with respect to the state-of-the-art semantic parsers on both datasets. On the ATIS dataset, our method outperforms STC, HVS and MLN parsers by 1.27%, 2.75%, and 5.44% respectively [2, 8, 11]. We also show that our method outperforms the handcrafted Phoenix parser by 3.28% on ASR output of the TownInfo dataset [2].

Although the TBL approach does not directly offer an N-best list of hypotheses with confidece scores, several methods have been developed to alleviate this problem. For example, transformation rules have been converted into decision trees [14] from which informative probability distributions on the class labels can be obtained. In future work, we will investigate how to obtain multiple hypotheses and confidence scores.

## 6. References

- [1] R. Kate, Y. Wong, and R. Mooney, "Learning to transform natural to formal languages," in *Proceedings of AAAI*, 2005.
- [2] F. Mairesse, M. M. Gašić, F. Jurčiček, S. Keizer, B. Thomson, K. Yu, and S. Young, "Spoken language understanding from unaligned data using discriminative classification models," in *Proceedings of ICASSP*, 2009.
- [3] D. Dahl, M. Bates, M. Brown, W. Fisher, K. Hunicke-Smith, D. Pallett, C. Pao, A. Rudnick, and E. Shriberg, "Expanding the scope of the ATIS task: The atis-3 corpus," in *Proceedings of the ARPA HLT Workshop*, 1994.
- [4] J. Williams and S. Young, "Partially observable markov decision for spoken dialog systems," *Computer Speech and Language*, vol. 21, no. 2, pp. 231–422, 2007.
- [5] B. Thomson, M. Gašić, S. Keizer, F. Mairesse, J. Schatzmann, K. Yu, and S. Young, "User study of the Bayesian update of dialogue state approach to dialogue management," in *Proceedings of Interspeech*, 2008.
- [6] E. Brill, "Transformation-Based Error-Driven Learning and Natural Language Processing: A Case Study in Part-of-Speech Tagging," *Computational Linguistics*, vol. 21, no. 4, pp. 543–565, 1995.
- [7] E. Briscoe, J. Carroll, and R. Watson, "The second release of the RASP system," in *Proceedings of COLING/ACL*, 2006.
- [8] Y. He and S. Young, "Semantic processing using the Hidden Vector State model," *Computer Speech & Language*, vol. 19, no. 1, pp. 85–106, 2005.
- [9] F. Jurčiček, J. Svec, and L. Muller, "Extension of HVS semantic parser by allowing left-right branching," in *Proceedings of ICASSP*, 2008.
- [10] L. Zettlemoyer and M. Collins, "Online learning of relaxed ccg grammars for parsing to logical form," in *Proceedings of EMNLP-CoNLL*, 2005.
- [11] I. Meza-Ruiz, S. Riedel, and O. Lemon, "Spoken language understanding in dialogue systems, using a 2-layer markov logic network: improving semantic accuracy," in *Proceedings of Londial*, 2008.
- [12] R. Kate, "A dependency-based word subsequence kernel," in *Proceedings of EMNLP*, 2008.
- [13] W. Ward, "The phoenix system: Understanding spontaneous *Proceedings of ICASSP*, 1991.
- [14] R. Florian, J. C. Henderson, and G. Ngai, "Coaxing confidence from an old friend: Probabilistic classifications from transformation rule lists," in *Proceedings EMNLP*, 2000.