

Multi-Task Learning for Parsing the Alexa Meaning Representation Language

Vittorio Perera

Carnegie Mellon
University
vdperera@cs.cmu.edu

Tagyoung Chung

Amazon Inc.
tagyoung@amazon.com

Thomas Kollar

Amazon Inc.
kollart@amazon.com

Emma Strubell

University of
Massachusetts Amherst
strubell@cs.umass.edu

Abstract

The Alexa Meaning Representation Language (AMRL) is a compositional graph-based semantic representation that includes fine-grained types, properties, actions, and roles and can represent a wide variety of spoken language. AMRL increases the ability of virtual assistants to represent more complex requests, including logical and conditional statements as well as ones with nested clauses. Due to this representational capacity, the acquisition of large scale data resources is challenging, which limits the accuracy of resulting models. This paper has two primary contributions. The first contribution is a linearization of the AMRL parses that aligns it to a related task of spoken language understanding (SLU) and a deep neural network architecture that uses multi-task learning to predict AMRL fine-grained types, properties and intents. The second contribution is a deep neural network architecture that leverages embeddings from the large-scale data resources that are available for SLU. When combined, these contributions enable the training of accurate models of AMRL parsing, even in the presence of data sparsity. The proposed models, which use the linearized AMRL parse, multi-task learning, residual connections and embeddings from SLU, decrease the error rates in the prediction of the full AMRL parse by 3.56% absolute.

Introduction

As intelligent assistants become more open and connected, there is a need to expand their capacity to understand task-oriented language. Expanding this capacity requires a representation that can handle many forms of spoken language and models that can predict this representation. This paper uses a new representation that can support a wide variety of spoken language called the Alexa Meaning Representation Language (AMRL). AMRL is a compositional, graph-based semantic representation that is backed by a large-scale ontology. An AMRL parse contains actions, fine-grained types, properties, and verb roles (as in Figure 1). AMRL provides the representational capacity to not only understand simple requests to a virtual assistant, but also complex logical, conditional, and nested statements.

An alternative representation to AMRL are those used for spoken language understanding (SLU), which categorize requests into domains, intents and slots (Figure 2). A domain

Copyright © 2018, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

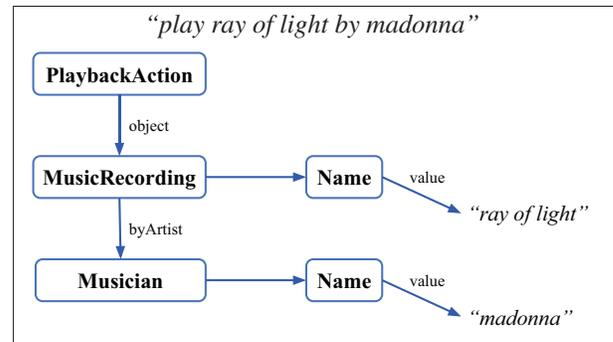


Figure 1: Meaning representation of a sentence using AMRL.

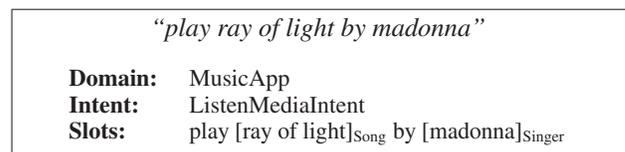


Figure 2: Example of the spoken language understanding (SLU) representation.

is a general category for a request (e.g., music, calendar), an intent is an action within that domain (e.g., play, search) and slots are mentions within the utterance (e.g., “ray of light” is the song to be played). By factoring the problem this way, each SLU domain may have its own unique set of intents and slots. In contrast, AMRL provides a common semantic representation for spoken language (Figure 3). There are three primary benefits to this common representation. The first benefit is that intents and slots that have the same meaning have the same labels. For example, “order me an echo dot” and “Alexa, order me a taxi” fall into distinct domains (e.g., shopping and taxi), resulting in distinct intents for each domain (e.g., OrderProduct and OrderTaxi). The fact that requests such as these, which have a similar surface form, can appear in different domains also makes it challenging to add new features without degrading SLU accuracy. The second benefit of AMRL is that cross-domain utterances are directly supported. For example, “find me a restaurant near the shark game” would belong to both the SLU local search and the

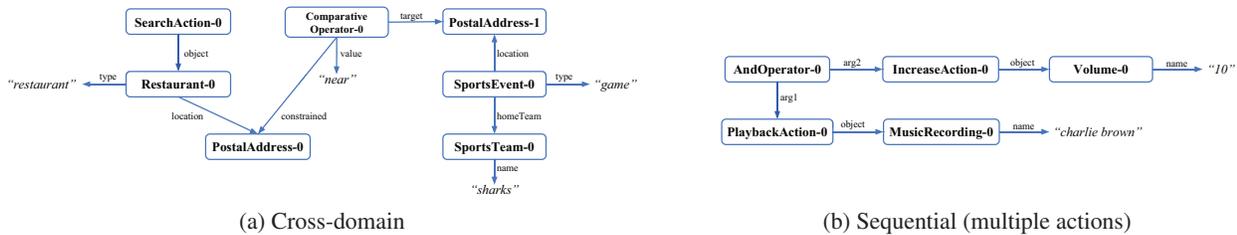


Figure 3: (a) is the AMRL for “find restaurants near the sharks game.”. (b) is the AMRL for “play charlie brown and turn the volume up to 10”. There are two actions that AMRL could handle with a sequential operator. These requests are challenging to handle using the current domain-based SLU representation.

<i>“play ray of light by madonna”</i>	
Properties:	play [ray of light]object.name by [madonna]object.byArtist.name
Entities:	play [ray of light]name@MusicRecording by [madonna]name@Musician
Intents:	PlaybackAction object@MusicRecording

Figure 4: Linearized annotation for AMRL.

sports domains. The third benefit is that more complex requests can be supported. For example, “play hunger games and turn the lights down to 3” is not easily supported existing SLU representation since it involves multiple domains and must also support sequential actions.

This paper presents models that can predict AMRL given a natural language utterance. In order to address the limited availability of AMRL data, a simple linearization scheme has been developed (Figure 4). The linearized AMRL format factors into three components: intents, types, and properties, which are predicted using a deep neural network architecture that is trained using multi-task learning. Each layer of the model predicts a component of the linearized AMRL representation. By ordering these tasks from coarse to fine, each subsequent layer is able to reuse representations from previous layers. The simplest model consists of a word-embedding layer and three bi-directional LSTM layers, residual connections, and a custom decoder to improve accuracy (He et al. 2017).

In addition to linearizing the AMRL parse graphs, we also leverage embeddings trained from the much larger datasets that are available for SLU. As a comparison, the linearized AMRL corpus consisted of 300k examples, while the SLU corpus consisted of around 3 million examples. To leverage these datasets, new layers for intents and slots were added to our DNN models, and were trained as new tasks using multi-task learning. Slot embeddings are found to produce a 2.6% improvement in IRER (full-parse accuracy), while domain embeddings improve the accuracy by another 0.1% IRER. An additional 1% IRER improvement is obtained by using a custom decoder that leverages span-based IOB and ontological constraints. Our proposed model even decreases the full-parse error rate by 1.5% IRER when compared to a

baseline model that has access to additional inputs, such as gazetteers and general-purpose word embeddings. The rest of this document is organized as follows, we review the related work, introduce the AMRL and SLU representations, the baseline and proposed models, a custom decoder and results.

Related Work

The Alexa ontology is a version of schema.org (Guha, Brickley, and Macbeth 2016) that has been adapted for spoken language understanding. The Alexa Ontology is used as a basis for the Alexa Meaning Representation Language. Some alternative semantic representations include FrameNet (Baker, Fillmore, and Lowe 1998), which is a semantic representation that represents an utterance, along with verb roles. Other approaches have a representation similar to AMRL including lambda-DCS (Liang 2013) and combinatory categorial grammars (CCG) (Steedman and Baldrige 2011). AMR is a representation that has a hierarchical graph-based approach that is well-suited for longer texts (Banarescu and others 2013; Kevin Knight 2017). Other approaches focus more on syntax than semantics, such as universal dependencies (Nivre and others 2016). Unlike these approaches, AMRL focuses on directly supporting spoken language understanding and contains fine-grained types along with actions, verb roles, and properties. An overview of other semantic representations is covered in (Abend and Rappoport 2017).

DNNs are widely used for sequence labeling. (Shimaoka et al. 2016) perform fine-grained entity labeling using a neural attention model. (Dong et al. 2015) use a combination of NNs to embed words and entities for coarse-grained entity labeling. More recently, two types of network architectures have gained popularity. The first one is LSTM-CNNs (Chiu and Nichols 2015), which use a combination of word-level and CNN-extracted character-level features to augment the input to bi-LSTMs. The second one is LSTM-CRFs (Huang, Xu, and Yu 2015), which apply a CRF constraints to bi-LSTMs. Recently, (Ma and Hovy 2016) combined the two approaches to get the state of the art results on standard CONLL 2003 NER task.

LSTMs (Long Short Term Memory) (Hochreiter and Schmidhuber 1997) perform well on many NLP tasks including sequence tagging, intent classification, and language modeling due to their inherent ability to model long term sequential dependencies. Bi-LSTMs (Graves, rahman Mo-

hamed, and Hinton 2013) are layered architectures which effectively use past and future information via Forward and Backward LSTM layers. Bi-LSTMs have been successfully applied to feature generation for tasks like dependency parsing (Kiperwasser and Goldberg 2016) and semantic role labeling (Zhou and Xu 2015). All our models adopt a deep neural-network architecture with Bi-LSTMs as our primary building blocks. For text classification, there has been a lot of recent interest in using character-level embeddings (Kim et al. 2015) as an additional input to neural architectures because of their ability to model morphological features as well as effectively handle out-of-vocabulary words. (Ballesteros, Dyer, and Smith 2015) use character embedding for dependency parsing. (Xiao and Cho 2016) combine character embeddings with CNNs for the text classification task.

Semantic parsing and spoken language understanding (SLU) learn to map natural language to a formal representation. Although semantic parsers can be trained using sentences annotated with this formal representation (Zelle and Mooney 1996; Zettlemoyer and Collins 2012; Wong and Mooney 2006; Kwiatkowski et al. 2010; Krishnamurthy and Mitchell 2012), they have not generally been used directly for spoken language. Most applications of spoken language understanding maps utterances onto a fixed domain, intent, and slot structure (Gupta et al. 2006). It cannot generally represent complex, cross-domain, or compositional utterances.

Multitask learning in deep neural networks has been shown to help generalization. Similar to our work but with CNNs (Xu and Sarikaya 2013) jointly model sentence classification and sequence labeling. (Guo et al. 2014) use recursive neural networks to jointly classify intents and fill slots. (Miwa and Bansal 2016) achieved state of the art for entity and relation classification. (Zhang and Weiss 2016) used them for part of speech tagging and dependency parsing. Transferring of learned embeddings was explored in in (Yosinski et al. 2014). Our work builds on these by creating a deep multi-task model for predicting a meaning representation for spoken language.

Approach

In this section we describe the Alexa meaning representation language (AMRL), our deep multi-task model, mechanisms for learning representations for spoken language understanding (SLU) and transferring those representations to improve AMRL parsing.

Alexa Meaning Representation Language

AMRL provides a common semantics for natural language understanding for voice applications. An AMRL parse consists of five primary components:

- **Actions** define the core functionality used for spoken language understanding. In Figure 1, the PlaybackAction is used on a MusicRecording, but can also be used on Books, Videos and other playable objects.
- **Roles** Actions operate on entities via roles. The .object role defines the entity on which the action operates.
- **Types** apply to each entity mention in a request and are hierarchical. In Figure 1, there is a MusicRecording and Musician type.

- **Properties** define relationships between variables of a given type. For example, the “byArtist” property of a MusicRecording, defines who wrote a song.
- **Operators** Operators can be used to represent complex logical or spatial relationships.

The aim of AMRL is to provide a common semantic representation for spoken language. It can represent anaphora, conditional statements, sequential actions, and logical expressions. AMRL can have arbitrary nesting, enabling it to represent complex statements. Figure 1 shows a simple example of AMRL. Figure 3 shows how sequential and cross-domain queries can be represented in AMRL.

In this paper, we investigate linearization as a means to address the data sparsity of annotated AMRL examples. Linearization is a common method to simplify syntactic (Vinyals et al. 2015), CCG (e.g., via supertagging) (Lewis, Lee, and Zettlemoyer 2016), and AMR parsing (Peng et al. 2017). The linearization addresses data sparsity in two primary ways. First, we use a linearization scheme similar to one that has been shown to improve parse accuracy for AMR parsing on similar size datasets (Peng et al. 2017). Secondly, we use a representation that aligns to the spoken language understanding problem, aligning these slots and intents with those used in AMRL parsing. This alignment should improve the alignment of embeddings from the SLU domain.

An example of the linearization scheme for AMRL can be seen in Figure 4. AMRL is a rooted graph. Starting at the root, the property linearization recursively descends to a leaf, appending each property visited along the way until a leaf node is reached. Incoming property arcs to a visited node are inverted to avoid cycles and handle multi-headed graphs. Types of the leaf nodes are the only ones that appear in the linearization. When a leaf node is visited, the corresponding property is included in the linearization (e.g., one of “type,” “value” or “name”). Intents include action, the roles on the action (e.g., “object”), and the type of each of the roles. This linearization enables the AMRL to be formulated as a sequential prediction problem, factoring into three components: intents, types, and properties as in Figure 4.

Models

AMRL parsing use multi-task learning to jointly learn tasks of predicting the SLU or AMRL representations. In this section we first describe the baseline model used in our experiments. Then we describe models that leverage representations from spoken language understanding. Compared to models trained separately for each of the tasks, multi-task models allow us to exploit commonalities between tasks (Caruana 1998) such as overlapping spans between AMRL types and properties.

Baseline model The baseline model is a deep bi-directional LSTM neural network trained using multi-task learning. There are three LSTM layers, each of which predicts a different component of the parse. The first layer performs type prediction, the next layer performs property prediction and the final layer performs intent classification. This structures the problem as a coarse-to-fine prediction problem. The types

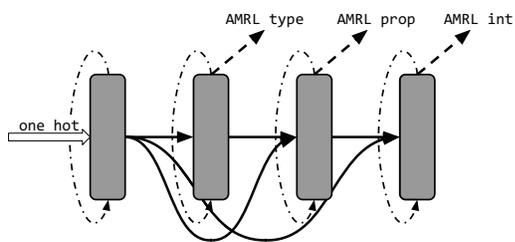


Figure 5: Topology of the baseline multi-task DNN model. Each component (above) is a bi-LSTM layer. The final layer is a softmax, and the dashed line is the recurrent connection. “type” and “prop” perform sequential prediction, while “int” categorizes the intent of the user.

are predicted at a coarse-grained level, after which the fine-grained properties and the actions are predicted. Residual connections are included to leverage embeddings from previous stages. For example, the word embeddings are used as input to the LSTM layers predicting the properties and the actions. Figure 5 shows the topology of this model.

Each block in Figure 5 is a bi-directional LSTM, which is used for both sequential prediction and classification. At the output layers, we form a prediction by feeding the concatenated hidden representations for the token we wish to label into another affine transform followed by a softmax to obtain posterior probabilities over the set of labels. Dropout is applied after each LSTM layer. For classification, only the final state from the forward and backward LSTM is used to predict the category.

The first layer provides a shared word embedding (Collobert and Weston 2008), while the remaining three LSTM layers are connected to an affine transform and a softmax layer to provide the output for each of the three tasks (i.e., properties, type, and intent prediction). IOB tagging is used to denote the inside, outside, and beginning of each property and type span. The input of this baseline model is a one-hot vector for each word in the sentence.

Transferring learned representations There are four primary modeling architectures that have been used to add embedding layers for domains, intents and slots. These architectures use domain and slot embeddings from the SLU task. Intent embeddings from the SLU task were added to initial models, but did not result significant performance gains.

The first two model architectures leverage embeddings from SLU slots. We investigate this first because we expect that SLU slots and AMRL types will often be correlated in terms of the spans and semantics. We expect that a shared representation that leverages this correlation is likely to improve the overall parse accuracy.

Two architectures that leverage a common embedding layer either in parallel or in serial to the AMRL type-prediction task. The parallel model (+SLU Slots), includes a separate SLU slot-prediction task in parallel to the AMRL type-prediction task resulting in a common word embedding layer. The resulting model can be seen in Figure 6a. The

serial model (+SLU Slots pipe) predicts the SLU slots before predicting the AMRL types, leveraging the embedding from the SLU task to directly predict the later stages of AMRL. This model can be seen in Figure 6b.

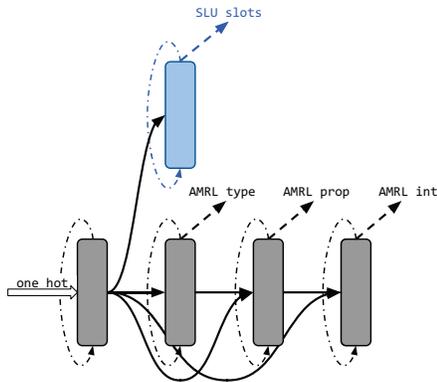
The remaining two modeling architectures leverage both domain and slot embeddings. The domain of a query is a strong indicator on the entities and their relationships. For example, if an utterance is from the “Music” domain, then it is likely include, for example, artists and songs, as well as certain relationships between the entities (e.g., “xxo” is by the artist “Coldplay”). As such, the remaining two architectures incorporate domain embeddings as input to the AMRL types and property layers. In both models, we only investigate serial embeddings for the SLU slots since they were found to be better than the parallel architecture in early experiments.

In the third model architecture (+SLU Slots/Domain), the domain layer is trained as a task in parallel to the AMRL type layer and is input to the property layer. The rationale behind this choice is that, although domain and property prediction are very different tasks, the domain, when correctly classified, can restrict what kind of properties we expect (e.g., if a sentence is classified in the *Music* domain we expect properties such as *byArtist* but not *weatherCondition*). There is still a shared embedding space learned as input to the type model. This model (+SLU Slots/Domain) is shown in Figure 6c.

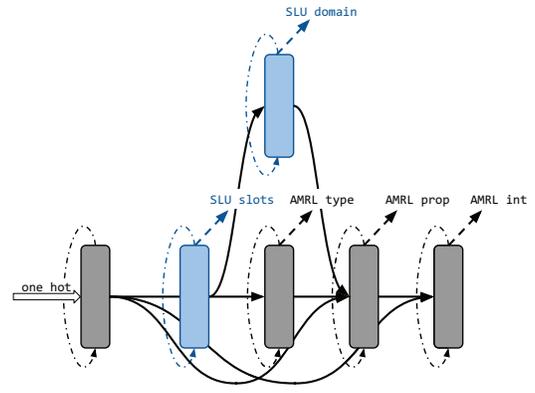
The final model architecture (+SLU Slots/Domain pipe), a common LSTM embedding layer is trained for SLU slots and domains. This layer is input to the layer that predicts the AMRL types. The hypothesis is that a single LSTM embedding layer might be enough to encode all the information from the SLU representation. Figure 6d shows the topology for this network.

Word embeddings and gazetteers We add to the input feature of our baseline pre-trained word embedding vectors and gazetteers. We use three hundred dimensional pre-trained word2vec embeddings, trained on the Google News corpus on 100 billion utterances (Mikolov et al. 2013) and we incorporate them as an additional input to the one-hot encoding of each word. Gazetteers (lists of entity mentions) from the Alexa ontology that backs AMRL were used in a similar way to word embeddings (e.g., as an additional input per word). For example, a Musician gazetteer will contain a list of music of musician names like “Sting.” These features are indicators that are set to 1 if the current word or word sequence appears in a gazetteer, and set to 0 otherwise. This model uses the same topology of the baseline, shown in Figure 5.

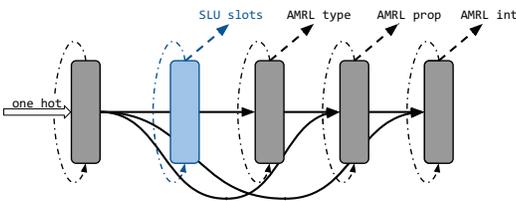
Decoding A beam-search decoder was written to leverage two primary constraints of this model in order to improve the accuracy of the result. The first constraint is to ensure the sequence of labels obey the IOB constraints. These constraints include (1) that an *I-* tag follows either another *I-* tag or a *B-* tag with the same label, and (2) a *B-* tag follows only an *O-* or an *I-* tag. The second constraint ensures that the final property label (e.g., name, type, or value) matches the initial property label, which ensures that we produce predictions consistent with the AMRL ontology. For example, in Figure 4, if “ray of light” were predicted as “type@MusicRecording”



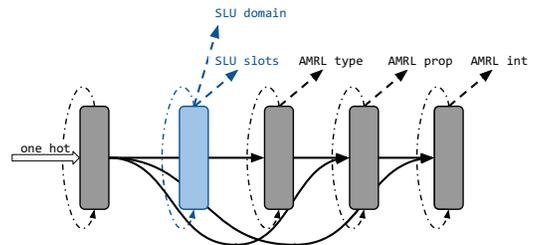
(a) Topology of the +SLU Slots model.



(c) Topology of the +SLU Slots/Domain model



(b) Topology of the +SLU Slots Pipe model



(d) Topology of the +SLU Slots/Domain Pipe model

Figure 6: Multi-task models used for training AMRL models alongside slot and domain embeddings from the existing SLU system. Highlighted in blue the bi-LSTM layers added compared to the baseline model.

but as a property was predicted at “object.name,” then this would be an invalid transition. The beam search limits property candidates to those that are above minimum probability. The value of the lower bound probability determines how aggressive the pruning is. The search is performed over the combination of properties and types created by combining the property candidates with all the possible matching types.

Optimization We optimized the parameters in two ways. In the first we fine-tuned pretrained embeddings for domain and slots. In the second, we performed joint training, learning the SLU embedding layers at the same time as the AMRL ones. Each model is trained until it fully converges on the training set, typically this takes around 60 epochs. We use a fixed learning rate of 0.0005 with an L2 penalty of $1e-8$, and a batch size of 128 sentences each. The output of each bi-LSTM layer is a vector of size 256 that is created by concatenating the hidden representation of forward and backward LSTMs (each of size 128).

To prevent overfitting to our training set we take two measures. First we connect the output of each bi-LSTM layer to a drop-out component with retention rate of 80%. The output of the drop-out component is then used as input for the following layers. Second we use weighted cross-entropy loss functions with a weight of λ when joint training (see below). We evaluate our accuracy on the development set to show

generalization performance.

Datasets

The two datasets used in these experiments are (1) a large corpus collected for spoken language understanding (SLU) and (2) a smaller corpus of linearized AMRL parses. The SLU corpus is composed by a total of $\sim 2.8m$ unique sentences. These sentences span on the order of 20 domains, 200 intents, and 200 slots. The vocabulary of this corpus amounts to $\sim 150k$ distinct words. The representation for this corpus is shown in Figure 2. The AMRL corpus is significantly smaller than the SLU corpus and contains only around 350k unique sentences in the linearized representation. This data consists of intents, types and properties. Figure 4 shows an example of the annotation for the same sentence shown in Figure 1. The AMRL corpus spans over 60 linearized intents, 200 properties and 200 types. The total vocabulary of this corpus is around 42k words, one order of magnitude smaller than the SLU one. The development set and test set contain around 48k sentences annotated using AMRL. Accuracy is reported only on the AMRL test set.

Metrics

Four metrics are considered to evaluate our models. The first two, $F1_{IC}$, $F1_{SC}$, are F1 scores evaluated respectively at

<i>Models</i>	$F1_{IC}$	$F1_{SC}$	<i>ICER</i>	<i>IRER</i>
Baseline	0.9383	0.8439	6.4464	25.7176
+SLU Slots	0.9416	0.8551	6.1254	24.2876
+SLU Slots (jt)	0.9389	0.8449	6.4587	25.6867
+SLU Slots pipe	0.9432	0.8602	5.8867	23.1312
+SLU Slots pipe (jt)	0.9390	0.8456	6.3311	25.3534
+SLU Slots/Domain	0.9431	0.8614	5.9649	23.0222
+SLU Slots/Domain (jt)	0.9435	0.8621	6.0204	23.1147
+SLU Slots/Domain pipe	0.9351	0.8232	6.8270	29.9418
+SLU Slots/Domain pipe (jt)	0.9400	0.8538	6.1789	24.4501

Table 1: Models and their results. Models marked as (jt) are trained using the joint training approach.

intent and slot level. $F1_{SC}$ is a strong metric, as it requires the spans and labels for both the property and the type tasks to be correct.

The Intent Classification Error Rate (ICER) is inversely correlated with the $F1_{IC}$ and measures the number of incorrect intents. ICER is not always sufficient as there may be multiple intents in an utterance. The formula for ICER is:

$$\frac{\#incorrect(intents)}{\#total\ intents}$$

Finally, the Intent Recognition Error rate (IRER) is computed as:

$$\frac{\#incorrect(interpretation)}{\#total\ utterances}$$

where we consider an interpretation incorrect if any of the slots or intents differs from the ground truth.

Results

Table 1 shows the results for different model architectures. The baseline model was trained using only the AMRL corpus. The vocabulary was pruned and only words appearing twice or more are used. Every other word is mapped to an “unknown” token resulting in a vocabulary of $\sim 25k$ words. From these results it appears that using SLU tasks clearly provides benefit when training AMRL models; our best model (+SLU Slots/Domain) outperforms the baseline across all the metrics and some by a considerable margin.

We also compared fine-tuning pretraining embeddings and joint training (learning the SLU embedding layers at the same time as the other embeddings). For pretraining, we train a network to predict the SLU tasks; once converged the AMRL-specific components (i.e., the last three LSTM layers of each model) are added and trained until convergence using a cross-entropy loss. Joint training optimizes both SLU tasks and AMRL ones at the same time. At each time-step a random training instance is selected from one of the two corpora with probability ($p = |AMRL|/|SLU|$). Since the size of the SLU corpus is much bigger than the AMRL one, we need to prevent overfitting on the SLU tasks. To do so we used weighted cross entropy loss functions where the SLU tasks were down-weighted by a factor $\lambda = |AMRL|/(10 \times |SLU|)$.

We see a slight improvement in accuracy of the models that use joint training, though there are no conclusive results.

For the first two models pretraining appears to result in higher accuracy but for the remaining two the opposite seems to hold true. One possible explanation for this behavior is the use of cross entropy loss function in conjunction with the joint training approach. In our experiments we fixed the weight for the loss function but additional hyper-parameter tuning might improve the overall result.

We also compared the accuracy of the baseline and our best performing model across different actions. Figure 7 shows that our dataset is skewed, with two of the actions (Playback and Search) covering more than 97% of the total training instances; the remaining 15 actions have an almost uniform amount of sentences. Table 2 shows the gain in ICER for each action. We observe how the improvement is modest on the two most represented actions but much more pronounced on the less represented ones. In general, we find that when there is sufficient data available (i.e., for the Playback and Search actions) adding more information from the SLU task is not very helpful. On the other hand, when fewer training instances are available the information provided by the SLU tasks becomes very valuable and strongly improves the results.

In Table 3, we evaluate the best models against those trained with word embeddings and gazetteers. The baseline

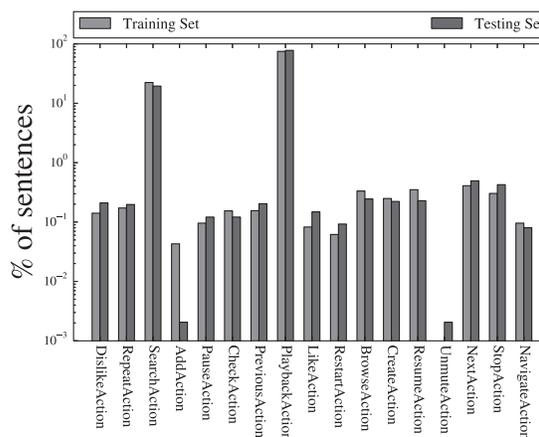


Figure 7: Action distribution on the training and testing set for the AMRL corpus.

Actions	Training instances	Δ ICER
PlaybackAction	HIGH	-0.53
SearchAction	MED	-0.17
NextAction	LOW	-3.34
ResumeAction	LOW	-10.81
BrowseAction	LOW	+5.88
StopAction	LOW	-1.93
CreateAction	LOW	0
RepeatAction	VLOW	+3.13
PreviousAction	VLOW	-23.23
CheckAction	VLOW	-5.08
DislikeAction	VLOW	-1.96
PauseAction	VLOW	-11.86
NavigateAction	VLOW	-5.12
LikeAction	VLOW	-6.94
RestartAction	VLOW	0

Table 2: Δ at ICER. The HIGH bin contains more than 100k examples. The MED bin contains between 2k and 100k examples. The LOW bin contains between 800 and 2k examples. The VLOW bin contains between 100 and 800 examples. The actions are ordered, in decreasing fashion, based on the number of occurrences in the AMRL training set.

model, which only has access only to the AMRL training dataset has lowest accuracy. Adding gazetteer and word embeddings improves accuracy, though the best model is the one that transfers the learned representations from the SLU task (proposed model from Table 1). Incorporating the additional constraints in the decoder (e.g., IOB and final property) results in our best model. For these experiments, we used a decoder with a beam of size 3 and a minimum probability of 10^{-7} . As expected the decoder does not impact any of the intent metrics ($F1_{IC}$ and ICER). The structurally incorrectness of the predicted outputs is upper bounded by 0.96% IRER using our proposed model with the custom decoder.

Conclusion

AMRL is a new graph-based representation for the meaning of a sentence. Since annotating AMRL is time consuming and costly, only a limited amount of data is available. In this paper we show that learned embeddings from related tasks can improve the accuracy of AMRL models. Domain and slot embeddings help significantly, improving the accuracy by 3.56% IRER (full-parse accuracy). A constrained decoder that leverages IOB and type/property constraints is a key component, decreasing IRER by 1% absolute. Although training time is a limiting factor, we would also like to explore the use of a CRF output layer to encode some of the decoding constraints.

References

Abend, O., and Rappoport, A. 2017. The state of the art in semantic representation. In *Proc. of ACL*.

Baker, C. F.; Fillmore, C. J.; and Lowe, J. B. 1998. The berkeley framenet project. In *Proceedings of the 36th Annual Meeting of the Association for Computational Linguistics*

	$F1_{IC}$	$F1_{SC}$	IRER
baseline	0.93	0.84	25.71
baseline+emd+gaz	0.94	0.86	23.65
proposed model	0.94	0.86	23.11
proposed model+decoder	0.94	0.87	22.15

Table 3: Results as compared to various baselines. Baseline is the multi-task model trained only on AMRL data. Baseline+emb+gaz is the baseline model with word embedding and gazetteer features as input to the model. Proposed model is the best model without word embeddings or gazetteer features. Proposed model+decoder includes results after decoding of the best model (without gazetteers or word embeddings as input). Beam size is three and a floor probability of 10^{-7} is used.

and 17th International Conference on Computational Linguistics - Volume 1, ACL '98, 86–90. Stroudsburg, PA, USA: Association for Computational Linguistics.

Ballesteros, M.; Dyer, C.; and Smith, N. A. 2015. Improved transition-based parsing by modeling characters instead of words with lstms. *arXiv preprint arXiv:1508.00657*.

Banarescu, L., et al. 2013. Abstract meaning representation for sembanking. In *Proceedings of the 7th Linguistic Annotation Workshop and Interoperability with Discourse*, 178–186. Sofia, Bulgaria: Association for Computational Linguistics.

Caruana, R. 1998. Multitask learning. In *Learning to learn*. Springer. 95–133.

Chiu, J. P., and Nichols, E. 2015. Named entity recognition with bidirectional lstm-cnns. *arXiv preprint arXiv:1511.08308*.

Collobert, R., and Weston, J. 2008. A unified architecture for natural language processing: Deep neural networks with multitask learning. In *Proceedings of the 25th international conference on Machine learning*, 160–167. ACM.

Dong, L.; Wei, F.; Sun, H.; Zhou, M.; and Xu, K. 2015. A hybrid neural model for type classification of entity mentions. In *Proceedings of the 24th International Conference on Artificial Intelligence*, 1243–1249. AAAI Press.

Graves, A.; rahman Mohamed, A.; and Hinton, G. 2013. Speech recognition with deep recurrent neural networks. In *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*.

Guha, R. V.; Brickley, D.; and Macbeth, S. 2016. Schema.org: Evolution of structured data on the web. *Commun. ACM* 59(2):44–51.

Guo, D.; Tur, G.; Yih, W.-t.; and Zweig, G. 2014. Joint semantic utterance classification and slot filling with recursive neural networks. In *Spoken Language Technology Workshop (SLT), 2014 IEEE*, 554–559. IEEE.

Gupta, N.; Tur, G.; Hakkani-Tur, D.; Bangalore, S.; Riccardi, G.; and Gilbert, M. 2006. The at&t spoken language understanding system. *IEEE Transactions on Audio, Speech, and Language Processing* 14(1):213–222.

He, L.; Lee, K.; Lewis, M.; and Zettlemoyer, L. 2017. Deep

- semantic role labeling: What works and whats next. In *Proceedings of the Annual Meeting of the Association for Computational Linguistics*.
- Hochreiter, S., and Schmidhuber, J. 1997. Long short-term memory. *Neural computation* 9(8):1735–1780.
- Huang, Z.; Xu, W.; and Yu, K. 2015. Bidirectional lstm-crf models for sequence tagging. In *arXiv preprint arXiv:1508.01991*.
- Kevin Knight, Bianca Badarau, L. B. C. B. M. B. K. G. U. H. D. M. M. P. T. O. N. S. 2017. Abstract meaning representation (amr) annotation release 2.0.
- Kim, Y.; Jernite, Y.; Sontag, D.; and Rush, A. M. 2015. Character-aware neural language models. *arXiv preprint arXiv:1508.06615*.
- Kiperwasser, E., and Goldberg, Y. 2016. Simple and accurate dependency parsing using bidirectional lstm feature representations. *Transactions of the Association for Computational Linguistics (TACL)* (4):313327.
- Krishnamurthy, J., and Mitchell, T. M. 2012. Weakly supervised training of semantic parsers. In *Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning*, 754–765. Association for Computational Linguistics.
- Kwiatkowski, T.; Zettlemoyer, L.; Goldwater, S.; and Steedman, M. 2010. Inducing probabilistic ccg grammars from logical form with higher-order unification. In *Proceedings of the 2010 Conference on Empirical Methods in Natural Language Processing, EMNLP '10*, 1223–1233. Stroudsburg, PA, USA: Association for Computational Linguistics.
- Lewis, M.; Lee, K.; and Zettlemoyer, L. 2016. Lstm ccg parsing. In *HLT-NAACL*, 221–231.
- Liang, P. 2013. Lambda dependency-based compositional semantics. *CoRR* abs/1309.4408.
- Ma, X., and Hovy, E. 2016. End-to-end sequence labeling via bi-directional lstm-cnns-crf. *arXiv preprint arXiv:1603.01354*.
- Mikolov, T.; Chen, K.; Corrado, G.; and Dean, J. 2013. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*.
- Miwa, M., and Bansal, M. 2016. End-to-end relation extraction using lstms on sequences and tree structures. *arXiv preprint arXiv:1601.00770*.
- Nivre, J., et al. 2016. Universal dependencies v1: A multilingual treebank collection. In Chair), N. C. C.; Choukri, K.; Declerck, T.; Goggi, S.; Grobelnik, M.; Maegaard, B.; Mariani, J.; Mazo, H.; Moreno, A.; Odijk, J.; and Piperidis, S., eds., *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC 2016)*. Paris, France: European Language Resources Association (ELRA).
- Peng, X.; Wang, C.; Gildea, D.; and Xue, N. 2017. Addressing the data sparsity issue in neural amr parsing. *arXiv preprint arXiv:1702.05053*.
- Shimaoka, S.; Stenetorp, P.; Inui, K.; and Riedel, S. 2016. Neural architectures for fine-grained entity type classification. *arXiv preprint arXiv:1606.01341*.
- Steedman, M., and Baldridge, J. 2011. *Combinatory Categorical Grammar*. Wiley-Blackwell. 181–224.
- Vinyals, O.; Kaiser, Ł.; Koo, T.; Petrov, S.; Sutskever, I.; and Hinton, G. 2015. Grammar as a foreign language. In *Advances in Neural Information Processing Systems*, 2773–2781.
- Wong, Y. W., and Mooney, R. J. 2006. Learning for semantic parsing with statistical machine translation. In *Proceedings of the main conference on Human Language Technology Conference of the North American Chapter of the Association of Computational Linguistics*, 439–446. Association for Computational Linguistics.
- Xiao, Y., and Cho, K. 2016. Efficient character-level document classification by combining convolution and recurrent layers. *arXiv preprint arXiv:1602.00367*.
- Xu, P., and Sarikaya, R. 2013. Convolutional neural network based triangular crf for joint intent detection and slot filling. In *Automatic Speech Recognition and Understanding (ASRU), 2013 IEEE Workshop on*, 78–83. IEEE.
- Yosinski, J.; Clune, J.; Bengio, Y.; and Lipson, H. 2014. How transferable are features in deep neural networks? In *Advances in neural information processing systems*, 3320–3328.
- Zelle, J. M., and Mooney, R. J. 1996. Learning to parse database queries using inductive logic programming. In *Proceedings of the national conference on artificial intelligence*, 1050–1055.
- Zettlemoyer, L. S., and Collins, M. 2012. Learning to map sentences to logical form: Structured classification with probabilistic categorical grammars. *arXiv preprint arXiv:1207.1420*.
- Zhang, Y., and Weiss, D. 2016. Stack-propagation: Improved representation learning for syntax. *arXiv preprint arXiv:1603.06598*.
- Zhou, J., and Xu, W. 2015. End-to-end learning of semantic role labeling using recurrent neural networks. In *Transactions of the Association for Computational Linguistics (TACL)*.