

# A Knowledge-Grounded Neural Conversation Model

Marjan Ghazvininejad,<sup>1\*</sup> Chris Brockett,<sup>2</sup> Ming-Wei Chang,<sup>2†</sup>  
Bill Dolan,<sup>2</sup> Jianfeng Gao,<sup>2</sup> Wen-tau Yih,<sup>2‡</sup> Michel Galley<sup>2</sup>

<sup>1</sup>Information Sciences Institute, USC

<sup>2</sup>Microsoft

ghazvini@isi.edu, mgalley@microsoft.com

## Abstract

Neural network models are capable of generating extremely natural sounding conversational interactions. However, these models have been mostly applied to casual scenarios (e.g., as “chatbots”) and have yet to demonstrate they can serve in more useful conversational applications. This paper presents a novel, *fully data-driven*, and knowledge-grounded neural conversation model aimed at producing more contentful responses. We generalize the widely-used Sequence-to-Sequence (SEQ2SEQ) approach by conditioning responses on both conversation history and external “facts”, allowing the model to be versatile and applicable in an open-domain setting. Our approach yields significant improvements over a competitive SEQ2SEQ baseline. Human judges found that our outputs are significantly more informative.

## Introduction

Recent work has shown that conversational chatbot models can be trained in an end-to-end and completely data-driven fashion, without hand-coding (Ritter, Cherry, and Dolan 2011; Sordoni et al. 2015; Shang, Lu, and Li 2015; Vinyals and Le 2015; Serban et al. 2016, *inter alia*). However, fully data-driven systems still lack grounding in the real world and do not have access to external knowledge (textual or structured), which makes it challenging for such systems to respond substantively. Fig. 1 illustrates the difficulty: while an ideal response would directly reflect on the entities mentioned in the query (user input), neural models produce responses that, while conversationally appropriate, seldom include factual content. This contrasts with traditional dialog systems, which can readily inject entities and facts into responses, but often at the cost of significant hand-coding. Slot-filler dialog systems are hard put to come up with a natural sounding utterance like the second response in Fig. 1 in a manner that is generalizable and scalable.

The goal of this work is to benefit from the versatility and scalability of fully data-driven models, while simultaneously seeking to produce models that are usefully grounded in external knowledge, permitting them to be deployed in, for ex-

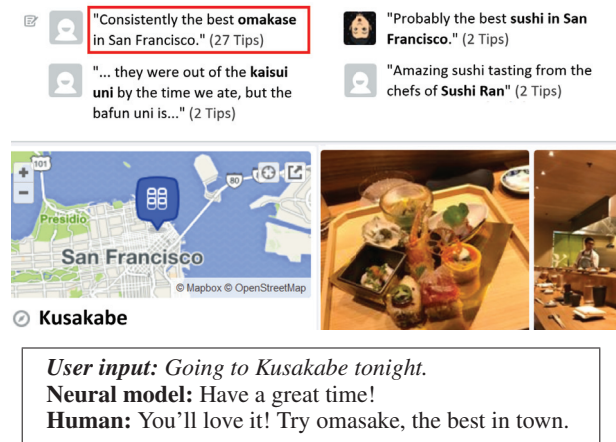


Figure 1: Responses of fully data-driven conversation models are often appropriate, but generally lack content characteristic of human responses.

ample, recommendation systems (e.g., for restaurants), and to adapt quickly and easily to new domains. The objective here is not task completion as in traditional dialog systems, but the ability to engage a user in a relevant and informative conversation. The tie to external data is critical, as the requisite knowledge is often not stored in conversational corpora. Much of this information is not found in structured databases either, but is textual, and can be mined from online resources such as Wikipedia, book reviews on Goodreads, and restaurant reviews on Foursquare.

This paper presents a novel, *fully data-driven*, knowledge-grounded neural conversation model aimed at producing contentful responses. Our framework generalizes the Sequence-to-Sequence (SEQ2SEQ) approach (Hochreiter and Schmidhuber 1997; Sutskever, Vinyals, and Le 2014) of previous neural conversation models, as it naturally combines conversational and non-conversational data via techniques such as multi-task learning (Caruana 1997; Liu et al. 2015). The key idea is that we can condition responses not only based on conversation history (Sordoni et al. 2015), but also on external “facts” that are relevant to the current context (for example, Foursquare entries as in Fig. 1). Our approach only requires a way to infuse external information

---

|    |  |
|----|--|
| A: | <b>Looking forward to trying @pizzalibretto tonight! my expectations are high.</b> |
| B: | <b>Get the rocco salad. Can you eat calamari?</b>                                  |

---

|    |   |
|----|---|
| A: | <b>Anyone in Chi have a dentist office they recommend? I'm never going back to [...] and would love a reco!</b> |
| B: | <b>Really looved Ora in Wicker Park.</b>  |

---

|    |  |
|----|--|
| A: | <b>I'm at California Academy of Sciences</b>   |
| B: | <b>Make sure you catch the show at the Planetarium. Tickets are usually limited.</b> |

---

|    |  |
|----|--|
| A: | <b>I'm at New Wave Cafe.</b>   |
| B: | <b>Try to get to Dmitri's for dinner. Their pan fried scallops and shrimp scampi are to die for.</b> |

---

|    |   |
|----|---|
| A: | <b>I just bought: [...] 4.3-inch portable GPS navigator for my wife, shh, don't tell her.</b> |
| B: | <b>I heard this brand loses battery power.</b>  |

---

Figure 2: Social media datasets include many contentful and useful exchanges, e.g., here recommendation dialog excerpts extracted from real tweets. While previous models (e.g., SEQ2SEQ) succeed in learning the **backbone of conversations**, they have difficulty modeling and producing *contentful words* such as named entities, which are sparsely represented in conversation data. To help solve this issue, we rely on non-conversational texts, which represent such entities much more exhaustively.

based on conversation context (e.g., via simple entity name matching), which makes it highly versatile and applicable in an open-domain setting. Using this framework, we have trained systems at a large scale using 23M general-domain conversations from Twitter and 1.1M Foursquare tips, showing significant improvements in terms of informativeness (human evaluation) over a competitive large-scale SEQ2SEQ model baseline. To the best of our knowledge, this is the first large-scale, fully data-driven neural conversation model that effectively exploits external knowledge.

### Related Work

The present work situates itself within the data-driven paradigm of conversation generation, in which statistical and neural machine translation models are derived from conversational data (Ritter, Cherry, and Dolan 2011; Sordoni et al. 2015; Serban et al. 2016; Shang, Lu, and Li 2015; Vinyals and Le 2015; Li et al. 2016a). The introduction of contextual models by (Sordoni et al. 2015) was an important advance within this framework, and we extend their basic approach by injecting side information from textual data. Introduction of side information has been shown to be beneficial to machine translation (Hoang, Cohn, and Haffari 2016), as has also the incorporation of images into multi-modal translation (Huang et al. 2016; Delbrouck, Dupont, and Seddati 2017). Similarly, (He et al. 2017) employ a knowledge graph to embed side information into dialog systems. Multi-task learning can be helpful in tasks ranging from query classification to machine translation (Caruana 1997; Dong et al. 2015; Liu et al. 2015; Luong et al. 2016). We adopt this approach in order to implicitly encode relevant external knowledge from textual data.

This work should be seen as distinct from more goal-directed neural dialog modeling in which question-answer slots are explicitly learned from small amounts of crowd-sourced data, customer support logs, or user data (Wen et al. 2015; 2017; Zhao et al. 2017). In many respects, that paradigm can be characterized as the neural extension of conventional dialog models with or without statistical modeling, e.g., (Oh and Rudnicki 2000; Ratnaparkhi 2002;

Banchs and Li 2012; Ameixa et al. 2014; Nio et al. 2014). Our purpose is to explore the space of less clearly goal-directed, but nonetheless informative (i.e., informational) conversation that does not demand explicit slot-filling.

Also relevant is (Bordes and Weston 2017), who employ memory networks to handle restaurant reservations, using a small number of keywords to handle entity types in a structured knowledge base. Similarly (Liu and Perez 2017) use memory networks to manage dialog state.

These works utilize datasets that are relatively small, and unlikely to scale, whereas we leverage free-form text to draw on datasets that are several orders of magnitude larger, allowing us to cover a greater diversity of domains and forms and thereby learn a more robust conversational backbone.

### Grounded Response Generation

A primary challenge in building fully data-driven conversation models is that most of the world's knowledge is not represented in any existing conversational datasets. While these datasets (Serban et al. 2015) have grown dramatically in size thanks in particular to social media (Ritter, Cherry, and Dolan 2011), such datasets are still very far from containing discussions of every entry in Wikipedia, Foursquare, Goodreads, or IMDB. This problem considerably limits the appeal of existing data-driven conversation models, as they are bound to respond evasively or deflectively as in Fig. 1, especially with regard to those entities that are poorly represented in the conversational training data. On the other hand, even where conversational data representing most entities of interest may exist, we would still face challenges as such huge dataset would be difficult to apply in model training, and many conversational patterns exhibited in the data (e.g., for similar entities) would be redundant.

Our approach aims to avoid redundancy and attempts to better generalize from existing conversational data, as illustrated in Fig. 2. While the conversations in the figure are about specific venues, products, and services, conversational patterns are general and equally applicable to other entities. The learned conversational behaviors could be used to, e.g., recommend other products and services. A traditional dia-

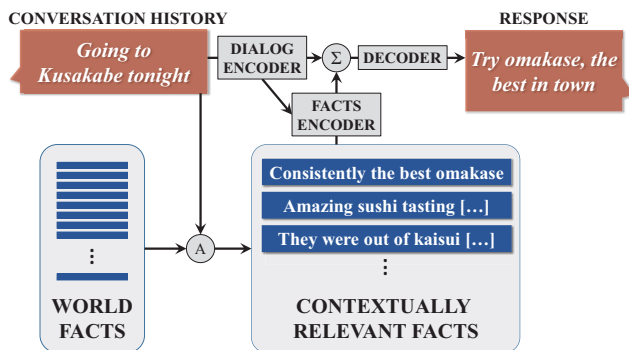


Figure 3: Knowledge-grounded model architecture.

log system would use predefined slots to fill conversational backbone (bold text) with content; here, we present a more robust and scalable approach.

In order to infuse the response with factual information relevant to the conversational context, we propose the knowledge-grounded model architecture depicted in Fig. 3. First, we have available a large collection of world facts,<sup>1</sup> which is a large collection of raw text entries (e.g., Foursquare, Wikipedia, or Amazon reviews) indexed by named entities as keys. Then, given a conversational history or source sequence  $S$ , we identify the “focus” in  $S$ , which is the text span (one or more entities) based on which we form a query to link to the facts. This focus can either be identified using keyword matching (e.g., a venue, city, or product name), or detected using more advanced methods such as entity linking or named entity recognition. The query is then used to retrieve all contextually relevant facts:  $F = \{f_1, \dots, f_k\}$ .<sup>2</sup> Finally, both conversation history and relevant facts are fed into a neural architecture that features distinct encoders for conversation history and facts. We will detail this architecture in the subsections below.

This knowledge-grounded approach is more general than SEQ2SEQ response generation, as it avoids the need to learn the same conversational pattern for each distinct entity that we care about. In fact, even if a given entity (e.g., @pizzalibretto in Fig. 2) is not part of our conversational training data and is therefore out-of-vocabulary, our approach is still able to rely on retrieved facts to generate an appropriate response. This also implies that we can enrich our system with new facts without the need to retrain the full system.

We train our system using multi-task learning (Caruana 1997; Luong et al. 2016) as a way of combining conversational data that is naturally associated with external data (e.g., discussions about restaurants and other businesses as in Fig. 2), and less informal exchanges (e.g., a response to *hi, how are you*). More specifically, our multi-task setup contains two types of tasks:

<sup>1</sup>For presentation purposes, we refer to these items as “facts”, but a “fact” here is simply any snippet of authored text, which may contain subjective or inaccurate information.

<sup>2</sup>In our work, we use a simple keyword-based IR engine to retrieve relevant facts from the full collection (see Datasets section).

- (1) one purely conversational, where we expose the model without fact encoder to  $(S, R)$  training examples,  $S$  representing the conversation history and  $R$  the response;
- (2) the other task exposes the full model with  $(\{f_1, \dots, f_k, S\}, R)$  training examples.

This decoupling of the two training conditions offers several advantages, including: First, it allows us to pre-train the conversation-only dataset separately, and start multi-task training (warm start) with a dialog encoder and decoder that already learned the backbone of conversations. Second, it gives us the flexibility to expose different kinds of conversational data in the two tasks. Finally, one interesting option is to replace the response in task (2) with one of the facts ( $R = f_i$ ), which makes task (2) similar to an autoencoder and helps produce responses that are even more contentful.

## Dialog Encoder and Decoder

The dialog encoder and response decoder form together a sequence-to-sequence (SEQ2SEQ model (Hochreiter and Schmidhuber 1997; Sutskever, Vinyals, and Le 2014), which has been successfully used in building end-to-end conversational systems (Sordani et al. 2015; Vinyals and Le 2015; Li et al. 2016a). Both encoder and decoder are recurrent neural network (RNN) models: an RNN that encodes a variable-length input string into a fixed-length vector representation and an RNN that decodes the vector representation into a variable-length output string. This part of our model is almost identical to prior conversational SEQ2SEQ models, except that we use gated recurrent units (GRU) (Chung et al. 2014) instead of LSTM (Hochreiter and Schmidhuber 1997) cells. Encoders and decoders in the present model do not share weights or word embeddings.

## Facts Encoder

The Facts Encoder of Fig. 3 is similar to the Memory Network model first proposed by (Weston, Chopra, and Bordes 2015; Sukhbaatar et al. 2015). It uses an associative memory for modeling the facts relevant to a particular problem—in our case, an entity mentioned in a conversation—then retrieves and weights these facts based on the user input and conversation history to generate an answer. Memory network models have been successfully used in Question Answering to make inferences based on the facts saved in the memory (Weston et al. 2016).

In our adaptation of memory networks, we use an RNN encoder to turn the input sequence (conversation history) into a vector, instead of a bag of words representation as used in the original memory network models. This enables us to better exploit interlexical dependencies between different parts of the input, and makes this memory network model (facts encoder) more directly comparable to a SEQ2SEQ model.

More formally, we are given an input sentence  $S = \{s_1, s_2, \dots, s_n\}$ , and a fact set  $F = \{f_1, f_2, \dots, f_k\}$  that are relevant to the conversation history. The RNN encoder reads the input string word by word and updates its hidden state. After reading the whole input sentence the hidden state of the RNN encoder,  $u$  is the summary of the input sentence.

By using an RNN encoder, we have a rich representation for a source sentence.

Let us assume  $u$  is a  $d$  dimensional vector and  $r_i$  is the bag of words representation of  $f_i$  with dimension  $v$ . Based on (Sukhbaatar et al. 2015) we have:

$$m_i = Ar_i \quad (1)$$

$$c_i = Cr_i \quad (2)$$

$$p_i = \text{softmax}(u^T m_i) \quad (3)$$

$$o = \sum_{i=1}^k p_i c_i \quad (4)$$

$$\hat{u} = o + u \quad (5)$$

Where  $A, C \in \mathbb{R}^{d \times v}$  are the parameters of the memory network. Then, unlike the original version of the memory network, we use an RNN decoder that is good for generating the response. The hidden state of the RNN is initialized with  $\hat{u}$  which is a symmetrization of input sentence and the external facts, to predict the response sentence  $R$  word by word.

As alternatives to summing up facts and dialog encodings in equation 5, we also experimented with other operations such as concatenation, but summation seemed to yield the best results. The memory network model of (Weston, Chopra, and Bordes 2015) can be defined as a multi-layer structure. In this task, however, 1-layer memory network was used, since multi-hop induction was not needed.

## Datasets

The approach we describe above is quite general, and is applicable to any dataset that allows us to map named entities to free-form text (e.g., Wikipedia, IMDB, TripAdvisor, etc.). For experimental purposes, we utilize datasets derived from two popular social media services: Twitter (conversational data) and Foursquare (non-conversational data).

**Foursquare:** Foursquare tips are comments left by customers about restaurants and other, usually commercial, establishments. A large proportion of these describe aspects of the establishment, and provide recommendations about what the customer enjoyed (or otherwise) We extracted from the web 1.1M tips relating to establishments in North America. This was achieved by identifying a set of 11 likely “foodie” cities and then collecting tip data associated with zipcodes near the city centers. While we targeted foodie cities, the dataset is very general and contains tips applicable to many types of local businesses (restaurants, theaters, museums, stores, etc.) In the interests of manageability for experimental purposes, we ignored establishments associated with fewer than 10 tips, but other experiments with up to 50 tips per venue yield comparable results. We limited the tips to those for which Twitter handles were found in the Twitter conversation data.

**Twitter:** We collected a **23M general dataset** of 3-turn conversations. This serves as a background dataset not associated with facts, and its massive size is key to learning the conversational structure or backbone.

Separately, on the basis of Twitter handles found in the Foursquare tip data, we collected approximately 1 million two-turn conversations that contain entities that tie to Foursquare. We refer to this as the **1M grounded dataset**. Specifically, we identify conversation pairs in which the first turn contained either a handle of the business name (preceded by the “@” symbol) or a hashtag that matched a handle.<sup>3</sup> Because we are interested in conversations among real users (as opposed to customer service agents), we removed conversations where the response was generated by a user with a handle found in the Foursquare data.

## Grounded Conversation Datasets

We augment the 1M grounded dataset with facts (here Foursquare tips) relevant to each conversation history. The number of contextually relevant tips for some handles can sometimes be enormous, up to 10k. To filter them for relevance to the input, the system vectorizes the input (as tf-idf weighted word counts) and each of the retrieved facts, and calculates cosine similarity between the input sentence and each of the tips and retains 10 tips with the highest score.

Furthermore, for a significant portion of the 1M Twitter conversations collected using handles found on Foursquare, the last turn was not particularly informative, e.g., when it provides a purely socializing response (e.g., *have fun there*). As one of our goals is to evaluate conversational systems on their ability to produce *contentful* responses, we select a dev and test set (4k conversations in total) designed to contain responses that are informative and useful.

For each handle, we created two scoring functions:

- Perplexity according to a 1-gram LM trained on all the tips containing that handle.
- $\chi$ -square score, which measures how much content each token bears in relation to the handle. Each tweet is then scored according to the average content score of its terms.

In this manner, we selected 15k top-ranked conversations using the LM score and 15k using the chi-square score. A further 15k conversations were randomly sampled. We then randomly sampled 10k conversations from these 45K conversations. Crowdsourced human judges were then presented with these 10K sampled conversations and asked to determine whether the response contained actionable information, i.e., did they contain information that would permit the respondents to decide, e.g., whether or not they should patronize an establishment. From this, we selected the top-ranked 4k conversations to be held out as validation set and test set; these were removed from our training data.

## Experimental Setup

### Multi-Task Learning

We use multi-task learning with these tasks:

- **FACTS** task: We expose the full model to  $(\{f_1, \dots, f_n, S\}, R)$  training examples.

<sup>3</sup>This mechanism of linking conversations to facts using exact match on the handle is high precision but low recall, but low recall seems reasonable as we are far from exhausting all available Twitter and Foursquare data.

- NOFACTS task: We expose the model without fact encoder to  $(S, R)$  examples.
- AUTOENCODER task: This is similar to the FACTS task, except that we replace the response with each of the facts, i.e., this model is trained on  $(\{f_1, \dots, f_n, S\}, f_i)$  examples. There are  $n$  times many samples for this task than for the FACTS task.<sup>4</sup>

The tasks FACTS and NOFACTS are representative of how our model is intended to work, but we found that the AUTOENCODER tasks helps inject more factual content into the response. The different variants of our multi-task learned system exploit these tasks as follows:

- SEQ2SEQ: Trained on task NOFACTS with the 23M general conversation dataset. Since there is only one task, it is not *per se* a multi-task setting.
- MTASK: Trained on two instances of the NOFACTS task, respectively with the 23M general dataset and 1M grounded dataset (but without the facts). While not an interesting system in itself, we include it to assess the effect of multi-task learning separately from facts.
- MTASK-R: Trained on the NOFACTS task with the 23M dataset, and the FACTS task with the 1M grounded dataset.
- MTASK-F: Trained on the NOFACTS task with the 23M dataset, and the AUTOENCODER task with the 1M dataset.
- MTASK-RF: Blends MTASK-F and MTASK-R, as it incorporates 3 tasks: NOFACTS with the 23M general dataset, FACTS with the 1M grounded dataset, and AUTOENCODER again with the 1M dataset.

We trained a one-layer memory network structure with two-layer SEQ2SEQ models. More specifically, we used 2-layer GRU models with 512 hidden cells for each layer for encoder and decoder, the dimensionality of word embeddings is set to 512, and the size of input/output memory representation is 1024. We used the Adam optimizer with a fixed learning rate of 0.1. Batch size is set to 128. All parameters are initialized from a uniform distribution in  $[-\sqrt{3/d}, \sqrt{3/d}]$ , where  $d$  is the dimension of the parameter. Gradients are clipped at 5 to avoid gradient explosion.

Encoder and decoder use different sets of parameters. The top 50k frequent types from conversation data is used as vocabulary which is shared between both conversation and non-conversation data. We use the same learning technique as (Luong et al. 2016) for multi-task learning. In each batch, all training data is sampled from one task only. For task  $i$  we define its mixing ratio value of  $\alpha_i$ , and for each batch we select randomly a new task  $i$  with probability of  $\alpha_i / \sum_j \alpha_j$  and train the system by its training data.

## Decoding and Reranking

We use a beam-search decoder similar to (Sutskever, Vinyals, and Le 2014) with beam size of 200, and maximum response length of 30. Following (Li et al. 2016a), we

<sup>4</sup>This is akin to an autoencoder as the fact  $f_i$  is represented both in the input and output, but is of course not strictly an autoencoder.

| Model     | Perplexity   |               |
|-----------|--------------|---------------|
|           | General Data | Grounded Data |
| SEQ2SEQ   | <b>55.0</b>  | 214.4         |
| SEQ2SEQ-S | 125.7        | <b>82.6</b>   |
| MTASK     | 57.2         | 82.5          |
| MTASK-R   | <b>55.1</b>  | <b>77.6</b>   |
| MTASK-F   | 77.3         | 448.8         |
| MTASK-RF  | 67.2         | 97.7          |

Table 1: Perplexity of different models. SEQ2SEQ-S is a SEQ2SEQ model that is trained on the NOFACTS task with 1M grounded dataset (without the facts).

generate  $N$ -best lists containing three features: (1) the log-likelihood  $\log P(R|S, F)$  according to the decoder; (2) word count; (3) the log-likelihood  $\log P(S|R)$  of the source given the response. The third feature is added to deal with the issue of generating commonplace and generic responses such as *I don't know*, which is discussed in detail in (Li et al. 2016a). Our models often do not need the third feature to be effective, but—since our baseline needs it to avoid commonplace responses—we include this feature in all systems. This yields the following reranking score:

$$\log P(R|S, F) + \lambda \log P(S|R) + \gamma |R|$$

$\lambda$  and  $\gamma$  are free parameters, which we tune on our development  $N$ -best lists using MERT (Och 2003) by optimizing BLEU. To estimate  $P(S|R)$  we train a Sequence-to-sequence model by swapping messages and responses. In this model we do not use any facts.

## Evaluation Metrics

Following (Sordani et al. 2015; Li et al. 2016a; Wen et al. 2017), we use BLEU automatic evaluation. While (Liu et al. 2016) suggest that BLEU correlates poorly with human judgment at the sentence-level,<sup>5</sup> we use instead corpus-level BLEU, which is known to better correlate with human judgments (Przybocki, Peterson, and Bronsart 2008), including for response generation (Galley et al. 2015). We also report perplexity and lexical diversity, the latter as a raw yet automatic measure of informativeness and diversity. Automatic evaluation is augmented with human judgments of appropriateness and informativeness.

## Results

**Automatic Evaluation:** We computed perplexity and BLEU (Papineni et al. 2002) for each system. These are shown in Tables 1 and 2 respectively. We notice that the SEQ2SEQ model specifically trained on general data has high perplexity on grounded data.<sup>6</sup> We observe that the perplexity of MTASK and MTASK-R models on both general

<sup>5</sup>This corroborates earlier findings that accurate sentence-level automatic evaluation is indeed difficult, even for Machine Translation (Graham, Baldwin, and Mathur 2015), as BLEU and related metrics were originally designed as corpus-level metrics.

<sup>6</sup>Training the system on just 1M grounded data with FACTS doesn't solve this problem, as its perplexity on general data is also high (not in table).

| Model    | BLEU        | Diversity    |              |
|----------|-------------|--------------|--------------|
|          |             | 1-gram       | 2-gram       |
| SEQ2SEQ  | 0.55        | 4.14%        | 14.4%        |
| MTASK    | 0.80        | 2.35%        | 5.9%         |
| MTASK-F  | 0.48        | 9.23%        | 26.6%        |
| MTASK-R  | <b>1.08</b> | 7.08%        | 21.9%        |
| MTASK-RF | 0.58        | <b>8.71%</b> | <b>26.0%</b> |

Table 2: BLEU-4 and lexical diversity.

and grounded data is as low as the SEQ2SEQ models that are trained specifically on general and grounded data respectively. As expected, injecting more factual content into the response in MTASK-F and MTASK-RF increased the perplexity especially on grounded data.

BLEU scores are low, but this is not untypical of conversational systems (e.g., (Li et al. 2016a; 2016b)). Table 2 shows that the MTASK-R model yields a significant performance boost, with a BLEU score increase of 96% and 71% jump in 1-gram diversity compared to the competitive SEQ2SEQ baseline. In terms of BLEU scores, MTASK-RF improvements is not significant, but it generates the highest 1-gram and 2-gram diversity among all models.

**Human Evaluation:** We crowdsourced human evaluations. We had annotators judge 500 randomly-interleaved paired conversations, asking them which was better on two parameters: appropriateness to the context, and informativeness. The crowd workers were instructed to: *Decide which response is more appropriate, i.e., which is the best conversational fit with what was said. Then decide which of the two is more informative (i.e., knowledgeable, helpful, specific) about the establishment under discussion.* Judges were asked to select among *Clearly #1, Maybe #Number 1, About the Same, Maybe #2, and Clearly #2*. These were converted to scores between 1 and 0, and assigned to the pair members depending on the order in which the pair was presented. Seven judges were assigned to each pair.<sup>7</sup>

The results of annotation are shown in Table 3. Our primary system MTASK-R, which performed best on BLEU, significantly outperforms the SEQ2SEQ baseline on *Informativeness* ( $p = 0.003$ ) and shows a small, but non-statistically-significant gain with respect *Appropriateness*. Other systems are included in the table for completeness. The “vanilla” MTASK shows no significant gain in *Informativeness*. MTASK-F performed significantly better than baseline ( $p = 0.005$ ) on *Informativeness*, but was significantly worse on *Appropriateness*. MTASK-RF came in slightly better than baseline on *Informativeness* but worse on *Appropriateness*, though in neither case is the difference statistically significant by the conventional standard of  $\alpha = 0.05$ . In sum, our best performing MTASK-R system appears to have successfully balanced the needs of informativeness and maintaining contextual appropriateness.

<sup>7</sup>Annotators whose variance fell greater than two standard deviations from the mean variance were dropped.

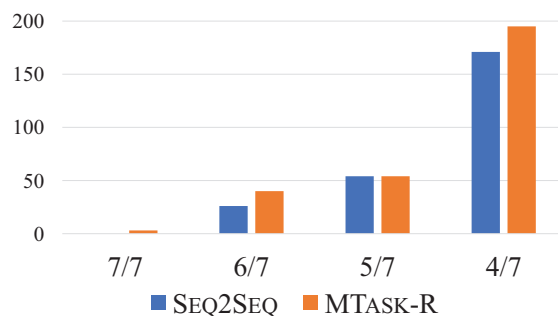


Figure 4: Judge preference counts (appropriateness) for MTASK-R versus SEQ2SEQ.

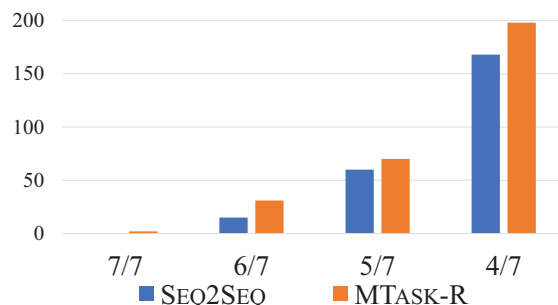


Figure 5: Judge preference counts (informativeness) for MTASK-R versus SEQ2SEQ.

The narrow differences in averages in Table 3 tend to obfuscate the judges’ voting trends. To clarify the picture, we translated the scores into the ratio of judges who preferred that system and binned the counts. Figs. 4 and 5 compare MTASK-R with the SEQ2SEQ baseline. Bin 7 on the left corresponds to the case where all 7 judges “voted” for the system, bin 6 to that where 6 out of 7 judges “voted” for the system, and so on.<sup>8</sup> Other bins are not shown since these are a mirror image of bins 7 through 4. The distributions in Fig. 5 are sharper and more distinctive than in Fig. 4, indicating that judge preference for the MTASK-R model is relatively stronger when it comes to informativeness.

## Discussion

Figure 6 presents examples from the MTASK-RF model, and illustrates that our responses are generally both appropriate and informative. First, we note that our models preserve the ability of earlier work (Sordani et al. 2015; Vinyals and Le 2015) to respond contextually and appropriately on a variety of topics, with responses such as *me too* (1) and *have a safe flight* (2). Second, our grounded models often incorporate information emerging from “facts”, while usually keeping the responses contextually appropriate. For example in (3), those facts revolve mostly around jewelry such as bracelets, which leads the system to respond *obsessed with their bracelets and the meaning behind them*,

<sup>8</sup>Partial scores were rounded up, affecting both systems equally.

| Comparison             | Informativeness |        |              |                | Appropriateness |                |       |        |
|------------------------|-----------------|--------|--------------|----------------|-----------------|----------------|-------|--------|
| SEQ2SEQ vs MTASK       | 0.501           | ±0.016 | 0.499        | ±0.015         | <b>0.530</b>    | ± <b>0.017</b> | 0.470 | ±0.017 |
| SEQ2SEQ vs MTASK-F     | 0.478           | ±0.015 | <b>0.522</b> | ± <b>0.015</b> | <b>0.537</b>    | ± <b>0.016</b> | 0.463 | ±0.017 |
| SEQ2SEQ vs MTASK-RF    | 0.492           | ±0.013 | 0.508        | ±0.013         | 0.502           | ±0.015         | 0.498 | ±0.014 |
| SEQ2SEQ vs MTASK-R (*) | 0.479           | ±0.017 | <b>0.521</b> | ± <b>0.013</b> | 0.495           | ±0.015         | 0.505 | ±0.015 |

Table 3: Mean differences in judgments in human evaluation, together with 95% confidence intervals. Differences sum to 1.0. Boldface items are significantly better ( $p < 0.05$ ) than their comparator. (\*): Main system, pre-selected on the basis of BLEU.

|    |  |
|----|--|
| 1  | A: I would give ALMOST anything for some [ <i>Mexican restaurant</i> ] right now.<br>B: Me too. Creamy jalapeno sauce is really good.  |
| 2  | A: Visiting the celebs at Los Angeles International Airport (LAX) - [...] w/ 70 others<br>B: Nice airport terminal. Have a safe flight.  |
| 3  | A: Obsessed with [ <i>jewelry company</i> ] :-*<br>B: oh my gosh obsessed with their bracelets and the meaning behind them !   |
| 4  | A: Live right now on [ <i>FM radio</i> ] Tune in!!!!<br>B: Listen to Lisa Paige  |
| 5  | A: Is that [...] in your photos? It's on my list of places to visit in NYC.<br>B: Don't forget to check out the 5th floor, while you are here, it's a great view.  |
| 6  | A: Are you [ <i>coworking community</i> ] today or tomorrow?<br>B: I think so. Try the cookies and coffee.   |
| 7  | A: Cant wait for my [ <i>pizzeria</i> ] tonight<br>B: Order bacon n pepperoni thin crust pizza.  |
| 8  | A: So [...] is down to one copy of Pound Foolish. I'm curious to see if they are re-ordering it.<br>B: Check out the video feed on 6 and take a picture of the Simpsons on the 3rd floor.                |
| 9  | A: I wish [...] would introduce another vegetarian option besides the shroomburger. It's delicious but kind of ridiculous.<br>B: This is the best j.crew in the world. Try the lemonade!                 |
| 10 | A: Just had an awesome dinner at [ <i>restaurant</i> ] Great recommendation [ <i>person</i> ]<br>B: One of my favorite places I've ever been to in NYC. The food is great and the service is lackluster. |

Figure 6: Responses produced by our neural response generation system MTASK-RF. Entities marked [...] are anonymized.

while *meaning behind them* does not belong to any of the facts and is instead “inferred” by the model (which associates jewelry with sentimental value). Responses influenced mostly by facts may occasionally contain a single unedited fact (4-5), but otherwise generally combine text from different facts (6-8).<sup>9</sup> Examples 9 and 10 are negative examples that illustrate the two main causes of system errors: the presence of an irrelevant fact (e.g., *j.crew* in example 9), and the system combining self-contradictory facts (10). Despite such errors, judges found that our best grounded system is generally on par with the SEQ2SEQ system in terms of appropriateness, while significantly improving informativeness (Table 3).

## Conclusions

We have presented a novel knowledge-grounded conversation engine that could serve as the core component of a

<sup>9</sup>Facts: *grab a cup of coffee and get productive and try the cookies in the vending machine of local food* (6); *sit with and take a picture of the Simpsons on the 3rd floor* and *Check out the video feed on 6 and Simpsons/billiards on 3!* (8).

multi-turn recommendation or conversational QA system. The model is a large-scale, scalable, fully data-driven neural conversation model that effectively exploits external knowledge, and does so without explicit slot filling. It generalizes the SEQ2SEQ approach to neural conversation models by naturally combining conversational and non-conversational data through multi-task learning. Our simple entity matching approach to grounding external information based on conversation context makes for a model that is informative, versatile and applicable in open-domain systems.

## Acknowledgements

We thank Xuetao Yin and Leon Xu for helping us obtain Foursquare data, and Kevin Knight, Chris Quirk, Nebojsa Jojic, Lucy Vanderwende, Vighnesh Shiv, Yi Luan, John Wieting, Alan Ritter, Donald Brinkman, and Puneet Agrawal.

## References

Ameixa, D.; Coheur, L.; Fialho, P.; and Quaresma, P. 2014. Luke, I am your father: dealing with out-of-domain requests by using movies subtitles. In *Intelligent Virtual Agents*.

- Banchs, R. E., and Li, H. 2012. IRIS: a chat-oriented dialogue system based on the vector space model. *ACL*.
- Bordes, A., and Weston, J. 2017. Learning end-to-end goal-oriented dialog. *ICLR 2017*.
- Caruana, R. 1997. Multitask learning. *Machine Learning* 28(1):41–75.
- Chung, J.; Gülgehre, Ç.; Cho, K.; and Bengio, Y. 2014. Empirical evaluation of gated recurrent neural networks on sequence modeling. *CoRR* abs/1412.3555.
- Delbrouck, J.-B.; Dupont, S.; and Seddati, O. 2017. Visually grounded word embeddings and richer visual features for improving multimodal neural machine translation. In *Grounding Language Understanding workshop*.
- Dong, D.; Wu, H.; He, W.; Yu, D.; and Wang, H. 2015. Multi-task learning for multiple language translation. *ACL*.
- Galley, M.; Brockett, C.; Sordoni, A.; Ji, Y.; Auli, M.; Quirk, C.; Mitchell, M.; Gao, J.; and Dolan, B. 2015. deltaBLEU: A discriminative metric for generation tasks with intrinsically diverse targets. *ACL-IJCNLP*.
- Graham, Y.; Baldwin, T.; and Mathur, N. 2015. Accurate evaluation of segment-level machine translation metrics. *NAACL*.
- He, H.; Balakrishnan, A.; Eric, M.; and Liang, P. 2017. Learning symmetric collaborative dialogue agents with dynamic knowledge graph embeddings. *ACL*.
- Hoang, C. D. V.; Cohn, T.; and Haffari, G. 2016. Incorporating side information into recurrent neural network language models. *NAACL-HLT*.
- Hochreiter, S., and Schmidhuber, J. 1997. Long short-term memory. *Neural computation* 9(8):1735–1780.
- Huang, P.-Y.; Liu, F.; Shiang, S.-R.; Oh, J.; and Dyer, C. 2016. Attention-based multimodal neural machine translation. *WMT*.
- Li, J.; Galley, M.; Brockett, C.; Gao, J.; and Dolan, B. 2016a. A diversity-promoting objective function for neural conversation models. *NAACL-HLT*.
- Li, J.; Galley, M.; Brockett, C.; Gao, J.; and Dolan, B. 2016b. A persona-based neural conversation model. *ACL*.
- Liu, F., and Perez, J. 2017. Dialog state tracking, a machine reading approach using memory network. *EACL*.
- Liu, X.; Gao, J.; He, X.; Deng, L.; Duh, K.; and Wang, Y.-Y. 2015. Representation learning using multi-task deep neural networks for semantic classification and information retrieval. *NAACL-HLT*.
- Liu, C.-W.; Lowe, R.; Serban, I.; Noseworthy, M.; Charlin, L.; and Pineau, J. 2016. How NOT to evaluate your dialogue system: An empirical study of unsupervised evaluation metrics for dialogue response generation. *EMNLP*.
- Luong, M.-T.; Le, Q. V.; Sutskever, I.; Vinyals, O.; and Kaiser, L. 2016. Multi-task sequence to sequence learning. *ICLR*.
- Nio, L.; Sakti, S.; Neubig, G.; Toda, T.; Adriani, M.; and Nakamura, S. 2014. Developing non-goal dialog system based on examples of drama television. In *Natural Interaction with Robots, Knowbots and Smartphones*. Springer.
- Och, F. J. 2003. Minimum error rate training in statistical machine translation. *ACL*.
- Oh, A. H., and Rudnicky, A. I. 2000. Stochastic language generation for spoken dialogue systems. *ANLP/NAACL Workshop on Conversational systems*.
- Papineni, K.; Roukos, S.; Ward, T.; and Zhu, W.-J. 2002. BLEU: a method for automatic evaluation of machine translation. *ACL*.
- Przybocki, M.; Peterson, K.; and Bronsart, S. 2008. Official results of the NIST 2008 metrics for machine translation challenge. In *MetricsMATR08 workshop*.
- Ratnaparkhi, A. 2002. Trainable approaches to surface natural language generation and their application to conversational dialog systems. *Computer Speech & Language* 16(3):435–455.
- Ritter, A.; Cherry, C.; and Dolan, W. B. 2011. Data-driven response generation in social media. *EMNLP*.
- Serban, I. V.; Lowe, R.; Charlin, L.; and Pineau, J. 2015. A survey of available corpora for building data-driven dialogue systems. *CoRR* abs/1512.05742.
- Serban, I. V.; Sordoni, A.; Bengio, Y.; Courville, A.; and Pineau, J. 2016. Building end-to-end dialogue systems using generative hierarchical neural network models. *AAAI*.
- Shang, L.; Lu, Z.; and Li, H. 2015. Neural responding machine for short-text conversation. *ACL-IJCNLP*.
- Sordoni, A.; Galley, M.; Auli, M.; Brockett, C.; Ji, Y.; Mitchell, M.; Nie, J.-Y.; Gao, J.; and Dolan, B. 2015. A neural network approach to context-sensitive generation of conversational responses. *NAACL-HLT*.
- Sukhbaatar, S.; Weston, J.; Fergus, R.; et al. 2015. End-to-end memory networks. *NIPS*.
- Sutskever, I.; Vinyals, O.; and Le, Q. V. 2014. Sequence to sequence learning with neural networks. *NIPS*.
- Vinyals, O., and Le, Q. 2015. A neural conversational model. *ICML*.
- Wen, T.-H.; Gasic, M.; Mrkšić, N.; Su, P.-H.; Vandyke, D.; and Young, S. 2015. Semantically conditioned LSTM-based natural language generation for spoken dialogue systems. *EMNLP*.
- Wen, T.-H.; Miao, Y.; Blunsom, P.; and Young, S. 2017. Latent intent dialog models. *ICML*.
- Weston, J.; Bordes, A.; Chopra, S.; Rush, A. M.; van Merriënboer, B.; Joulin, A.; and Mikolov, T. 2016. Towards AI-complete question answering: A set of prerequisite toy tasks. *ICLR*.
- Weston, J.; Chopra, S.; and Bordes, A. 2015. Memory networks. *ICLR*.
- Zhao, T.; Lu, A.; Lee, K.; and Eskenazi, M. 2017. Generative encoder-decoder models for task-oriented spoken dialog systems with chatting capability. *ACL*.