

Natural language is a fundamental form of information and communication. My ultimate research goal is to design Machine Learning approaches to Natural Language Processing systems that can understand language and communicate with people in different contexts, domains, text styles, and languages. To be specific, I focus on building intelligent systems that (1) first understand the contextual meaning of language grounded in various contexts where the language is used and (2) then generate effective language response in different forms for information request and human-computer communication. To this end, during my Ph.D. study, I decompose my long-term goal into the following research questions:

- How to enable machines to understand language at different levels of granularity including entities, sentences, documents, and conversations [1, 2, 3, 4]?
- How can agents summarize the information in the email, news, and scientific articles [5, 6, 7]?
- How can we synthesize formal programs (e.g., SQL) from natural language for question answering between humans and machines [8, 9, 10, 11, 12, 13, 14]?
- How to retrieve information relevant to the user query from documents in low-resource languages [15, 16]?

The challenges and solutions are inherently multi-disciplinary, spanning areas including Natural Language Processing, Deep Learning, Information Retrieval, Database, and Human-Computer Interaction. The sections below present my completed projects towards the answers to each question published as 16 papers in top AI conferences, followed by future directions in which I would like to make progress.

**Deep Neural Modeling of Text Units.** Natural language consists of text units at different granularity levels such as entity mentions, sentences, documents, and conversations. Small text units make up larger ones to convey meaning, and the semantics of textual units depend on each other. In particular, I am interested in how deep learning models can help us better understand and utilize the dependency relationship among text units in different NLP tasks. I have designed end-to-end deep neural networks for (1) entity extraction and coreference resolution in documents [2], (2) sentiment analysis and text classification for sentences and documents [1], and (3) addressee and response selection for multi-turn multi-party conversations [3]. The third one is described in detail as follows.

thread id	Sender	Addressee	Utterance
1	codepython	wafflejock	thanks
1	wafflejock	codepython	yup np
2	wafflejock	theoletom	you can use sudo apt-get install packagename – reinstall, to have apt-get install reinstall some package/metapackage and redo the configuration for the program.
3	codepython	-	i installed ubuntu on a separate external drive. now when i boot into mac, the external drive does not show up as bootable. the blue light is on. any ideas?
4	Guest54977	-	hi. wondering who knows where an ubuntu backup can be retrieved from.
2	theoletom	wafflejock	it's not a program. it's a desktop environment.
4	Guest54977	-	did some searching on my system and googling, but couldn't find an answer
2	theoletom	-	be a trace of it left yet there still is.
5	releaf	-	what's your opinion on a \$500 laptop that will be a dedicated ubuntu machine?
3	codepython	-	my usb stick shows up as bootable (efi) when i boot my mac. but not my external hard drive on which i just installed ubuntu. how do i make it bootable from mac hardware?
3	Jordan_U	codepython	did you install ubuntu to this external drive from a different machine?
5	Umeaboy	releaf	what country you from?
5	wafflejock		
Model Prediction		Addressee	Response
SI-RNN		★ releaf	★ there are a few ubuntu dedicated laptop providers like umeaboy is asking depends on where you are

Table 1: An example of addressee and response selection in Ubuntu IRC multi-party dialog. SI-RNN chooses to engage in a new sub-conversation by suggesting a solution to “releaf” about Ubuntu dedicated laptops. ★ denotes the ground-truth.

Real-world conversations often involve more than two speakers. In the Ubuntu Internet Relay Chat channel (Ubuntu IRC), for example, one user can initiate a discussion about an Ubuntu-related technical issue, and many other users can work together to solve the problem. Therefore, a multi-party dialog can have complex speaker interactions: at each turn, users play one of three roles (sender, addressee, observer), and those roles vary across turns. In this project, I study the problem of addressee and response selection in multi-party conversations: given a responding speaker and a dialog context, the task is to select an addressee and a response from a set of candidates for the responding speaker. Our model output for a task example is given in Table 1. The task requires modeling multi-party conversations and can be directly used to build retrieval-based dialog systems.

To model the complexity of multi-party dialogues, I designed the Speaker Interaction Recurrent Neural Network (SI-RNN). SI-RNN uses its dialog encoder to maintain speaker embeddings in a role-sensitive way. Speaker embeddings are updated in different GRU-based units based on their roles (sender, addressee, observer). Furthermore, noting that the addressee and response are mutually dependent, SI-RNN models the conditional probability (of addressee given the response and vice versa) and selects the addressee and response pair by maximizing the joint probability. On the public Ubuntu IRC benchmark data set, SI-RNN significantly improves the addressee and response selection performance by 10% accuracy, particularly in complex conversations with many speakers and responses to distant messages many turns in the past.

**Text Summarization in Different Domains and Styles.** Summarization aims to produce fluent and coherent synopses covering salient information in the documents. Recently, neural methods have shown promising results in text summarization using both extractive and abstractive approaches. However, most of the work focuses on the single-document setting in the news domain, relying on large training datasets such as the Gigaword Corpus and the CNN/Daily Mail.

My research has been focused on expanding neural summarization methods to various text domains and styles. I have worked on summarization in the following scenarios: (1) Use graph convolutional neural networks to summarize multiple long news articles [5] (co-lead with Michihiro Yasunaga), (2) Incorporate citation information for summarizing scientific articles [6] (led by Michihiro Yasunaga), and (3) Generate subject line for personal email messages [7] by optimizing quality estimation scores via reinforcement learning. The third project is described below.

Email is a ubiquitous form of communication. An email message consists of two elements: an *email subject line* and an *email body*. The subject line should tell what the email body is about and what the sender wants to convey. Table 2 shows an email body with three possible subject lines. I proposed the task of Subject Line Generation (SLG): automatically producing email subjects given the email body. Compared with news headline generation or news single document summarization, email subjects are generally much shorter, which means a system must have the ability to summarize with a high compression ratio.

To introduce the task, I built the first dataset, Annotated Enron Subject Line Corpus (AESLC). Furthermore, in order to properly evaluate the subject, I used a combination of automatic metrics from the text summarization and machine translation fields, in addition to building my own regression-based Email Subject Quality Estimator (ESQE). Third, to generate effective email subjects, I proposed a method that combines extractive and abstractive summarization using a two-stage process by Multi-Sentence Selection and Rewriting with Email Subject Quality Estimation Reward. The multi-sentence extractor first selects multiple sentences from the input email body. Extracted sentences capture salient information for writing a subject such as named entities and dates. Then, the multi-sentence abstractor rewrites multiple selected sentences into a succinct subject line while preserving key information. For training the network, I used a multi-stage training strategy incorporating both supervised cross-entropy training and reinforcement learning (RL) by optimizing the reward provided by the ESQE model. The automatic and human evaluations demonstrated that the model outperformed competitive baselines and approaches human-level quality.

**Multi-turn Text-to-SQL Semantic Parsing.** Generating SQL queries from user utterances is important to help people acquire information from databases. Such a text-to-SQL semantic parsing system bridges the data and the user through an intelligent natural language interface, greatly improving the efficiency of querying databases for many users beyond database experts. Furthermore, in a real-world application, users often access information in a multi-turn interaction with the system by asking a sequence of related questions. The users may explicitly refer to or omit previously mentioned entities and constraints, and may introduce refinements, additions or substitutions to what has already been said. This requires a text-to-SQL system to effectively process context information to synthesize the correct SQL.

To advance the research progress in this field, we built two multi-turn text-to-SQL data sets: (1) SPaC [9] (Figure 1, co-lead with Tao Yu) for cross-domain **S**emantic **P**arsing in **C**ontext. It contains 4,298 unique question sequences with 12k+ questions annotated with SQL queries. (2) CoSQL [10] (Figure 2, co-lead with Tao Yu) for building database-querying dialogue systems. It consists of 30k+ turns plus 10k+ annotated SQL queries, obtained from a Wizard-of-Oz collection of 3k dialogues. Both of them are built on top of our Spider dataset [11] (co-lead with Tao Yu), the largest cross-domain context-independent text-to-SQL dataset available in the field, and thus span 200 complex databases over

---

<b>Email Body:</b>	Hi All, I would be grateful if you could get to me today via email a job description for your current role. I would like to get this to the immigration attorneys so that they can finalise the paperwork in preparation for INS filing once the UBS deal is signed. Kind regards,
<b>Subject 1:</b>	Current Job Description Needed ( <i>COMMENT: This is good because it is both informative and succinct.</i> )
<b>Subject 2:</b>	Job Description ( <i>COMMENT: This is okay but not informative enough.</i> )
<b>Subject 3:</b>	Request ( <i>COMMENT: This is bad because it does not contain any specific information about the request.</i> )

---

Table 2: An email with three possible subject lines.

$D_1$  : Database about student dormitory containing 5 tables.

$C_1$  : Find the first and last names of the students who are living in the dorms that have a TV Lounge as an amenity.

$Q_1$  : How many dorms have a TV Lounge?

$S_1$  : SELECT COUNT(\*) FROM dorm AS T1 JOIN has\_amenity AS T2 ON T1.dormid = T2.dormid JOIN dorm\_amenity AS T3 ON T2.amenid = T3.amenid WHERE T3.amenity\_name = 'TV Lounge'

$Q_2$  : What is the total capacity of these dorms?

$S_2$  : SELECT SUM(T1.student\_capacity) FROM dorm AS T1 JOIN has\_amenity AS T2 ON T1.dormid = T2.dormid JOIN dorm\_amenity AS T3 ON T2.amenid = T3.amenid WHERE T3.amenity\_name = 'TV Lounge'

$Q_3$  : How many students are living there?

$S_3$  : SELECT COUNT(\*) FROM student AS T1 JOIN lives\_in AS T2 ON T1.stuid = T2.stuid WHERE T2.dormid IN (SELECT T3.dormid FROM has\_amenity AS T3 JOIN dorm\_amenity AS T4 ON T3.amenid = T4.amenid WHERE T4.amenity\_name = 'TV Lounge')

$Q_4$  : Please show their first and last names.

$S_4$  : SELECT T1.fname, T1.lname FROM student AS T1 JOIN lives\_in AS T2 ON T1.stuid = T2.stuid WHERE T2.dormid IN (SELECT T3.dormid FROM has\_amenity AS T3 JOIN dorm\_amenity AS T4 ON T3.amenid = T4.amenid WHERE T4.amenity\_name = 'TV Lounge')

Figure 1: Two question sequences from SPaC. Questions ( $Q_i$ ) in each sequence query a database ( $D_m$ ), obtaining information sufficient to complete the interaction goal ( $C_m$ ). Each question is annotated with a SQL query ( $S_i$ ). SQL segments from the interaction context are underlined.

138 domains. The large number of domains provide rich contextual phenomena and thematic relations between the questions, which general-purpose natural language interfaces to databases have to address. In addition, it enables us to test the generalization of the trained systems to unseen databases and domains. We are actively maintaining the datasets and leaderboards of our Text-to-SQL Challenge Series including Spider (<https://yale-lily.github.io/spider>), SPaC (<https://yale-lily.github.io/sparc>), and CoSQL (<https://yale-lily.github.io/cosql>).

Furthermore, I also proposed an editing-based approach for our cross-domain multi-turn text-to-SQL generation task [8]. Based on the observation that adjacent natural language questions are often linguistically dependent and their corresponding SQL queries tend to overlap, I utilized the interaction history by editing the previous predicted query to improve the generation quality. This editing mechanism views SQL as sequences and reuses generation results at the token level in a simple manner. It is flexible to change individual tokens and robust to error propagation. Furthermore, to deal with complex table structures in different domains, I employed an utterance-table encoder and a table-aware decoder to incorporate the context of the user utterance and the table schema. Experiment results on SPaC showed that by generating from the previous query, the model delivered an improvement of 7% question match accuracy and 11% interaction match accuracy over the previous state-of-the-art.

**Low-Resource Cross-lingual Information Retrieval.** Cross-Lingual Information Retrieval (CLIR) is the task of ranking foreign documents against a user query. As multilingual documents are more accessible, CLIR is increasingly more important whenever the relevant information is in other languages. Traditional CLIR systems consist of two components: machine translation and monolingual information retrieval. In this way, we first solve the translation problem, and the task is transformed to the monolingual setting. However, while conceptually simple, the performance of this modular approach is fundamentally limited by the quality of machine translation.

Recently, many deep neural learning-to-rank models have shown promising results in information retrieval. They learn a scoring function directly from the relevance label of query-document pairs. However, most previous neural IR papers only work with monolingual datasets, primarily for two reasons. First, when queries and documents are in different

$D_1$  : Database about student dormitories containing 5 tables

$Q_1$  : What are the names of all the dorms? INFORM\_SQL

$S_1$  : SELECT dorm\_name FROM dorm

$A_1$  : (Result table with many entries)

$R_1$  : This is the list of the names of all the dorms. CONFIRM\_SQL

$Q_2$  : Which of those dorms have a TV lounge? INFORM\_SQL

$S_2$  : SELECT T1.dorm\_name FROM dorm AS T1 JOIN has\_amenity AS T2 ON T1.dormid = T2.dormid JOIN dorm\_amenity AS T3 ON T2.amenid = T3.amenid WHERE T3.amenity\_name = 'TV Lounge'

$A_2$  : (Result table with many entries)

$R_2$  : This shows the names of dorms with TV lounges. CONFIRM\_SQL

$Q_3$  : What dorms have no study rooms as amenities? AMBIGUOUS

$R_3$  : Do you mean among those with TV Lounges? CLARIFY

$Q_4$  : Yes. AFFIRM

$S_4$  : SELECT T1.dorm\_name FROM dorm AS T1 JOIN has\_amenity AS T2 ON T1.dormid = T2.dormid JOIN dorm\_amenity AS T3 ON T2.amenid = T3.amenid WHERE T3.amenity\_name = 'TV Lounge' EXCEPT SELECT T1.dorm\_name FROM dorm AS T1 JOIN has\_amenity AS T2 ON T1.dormid = T2.dormid JOIN dorm\_amenity AS T3 ON T2.amenid = T3.amenid WHERE T3.amenity\_name = 'Study Room'

$A_4$  : Fawty Towers

$R_4$  : Fawty Towers is the name of the dorm that has a TV lounge but not a study room as an amenity. CONFIRM\_SQL

$Q_8$  : Thanks! THANK\_YOU

$R_8$  : You are welcome. WELCOME

Figure 2: A dialog from CoSQL. Gray boxes separate the user inputs ( $Q_i$ ) querying the database ( $D_i$ ) from the SQL queries ( $S_i$ ), returned answers ( $A_i$ ), and expert responses ( $R_i$ ).

languages, it is not clear how to measure the similarity of them in distributed representation space. Furthermore, deep neural networks need a large amount of training data to achieve decent performance. The annotation is prohibitively expensive for low-resource language pairs in our cross-lingual case.

In this paper [15], I proposed a cross-lingual deep relevance ranking architecture based on a bilingual view of queries and documents. As shown in Figure 3, the model first translates queries and documents and then uses four components to match them in both the source and target language. Each component is implemented as a term interaction network because they can make use of cross-lingual embeddings to explicitly encode terms of queries and documents in different languages. The final relevance score combines all components which are jointly trained given the relevance label. To deal with the small amount of training data, I first performed query likelihood retrieval and included the score as an extra feature in the model. In this way, the model effectively learns to rerank from only a few hundred relevance labels. Furthermore, by aligning word embedding spaces for multiple languages, the model can be directly applied under a zero-shot transfer setting when no training data is available for another language pair. On the MATERIAL CLIR dataset with three language pairs including English to Swahili, English to Tagalog, and English to Somali, the model outperformed other translation-based query likelihood retrieval and monolingual deep relevance ranking approaches by 2%-4% mAP scores.

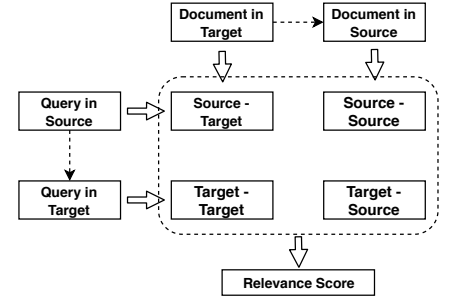


Figure 3: Cross-lingual Relevance Ranking with Bilingual Query and Document Representation.

## Future Research Directions

Building upon my past work, I plan to explore the following new directions and challenges.

**Grounding Language to User Interface Actions.** Mobile applications have been widely used in daily life. To accomplish a task, the user interacts with an application via a sequence of low-level user interface actions, such as click an element, swipe to check the rest of a list, and type a string. I am interested in building the agent that receives a high-level user instruction in natural language learns to perform a sequence of actions to accomplish the goal. This would automate the task completion that currently requires complicated human participation, greatly promoting the efficiency and accessibility for mobile application users. The complex representation of a mobile phone application contains both structured properties (e.g., the internal tree representation of a screen and the spatial relations among elements) and unstructured elements (e.g., text and image icons). Moreover, the agent needs to reason about the semantics of a high-level user goal and the context of screens and then learns to perform the correct sequence of actions. Therefore, this presents a challenging natural language grounding and navigation problem in this diverse and open-domain platform.

**Multilingual Transfer Learning for Low-resource Languages.** Deep neural networks still heavily rely on large amounts of in-domain labeled training data. However, only a few languages have large amounts of labeled data, and generalization in low-resource scenarios is still an open challenge. I would like to develop multilingual transfer learning techniques that can leverage annotations in high-resource languages to boost the performance of low-resource languages. I am particularly interested in models that require minimal cross-lingual supervision, leverage knowledge from multiple source languages, and quickly adapt to the target language and task.

**Natural Language for AI Interpretability.** While deep learning has become the de facto approach to build intelligent systems, the improvement of performance often comes at a cost of interpretability. Complex neural networks permit easy architectural and operational variations for state-of-the-art accuracy, yet they provide little transparency about their inner decision-making mechanisms. I am interested in how natural language can promote interpretable AI: language is not only the means of communication between humans, but it also offers a media for the intelligent system to explain and rationalize its solutions. To this end, I would like to empower the intelligent systems with abilities to automatically extract or generate human-readable language explanations to justify their predictions or actions.

**Controllable and Personalized Text Generation.** While deep learning models have shown promising results in text generation tasks such as summarization, translation, and dialog response generation, progress remains to be made towards controllable and personalized text generation. In particular, I would like to develop more controllable models such that (1) the generated text stay faithful to the conditioned text input, (2) we can manipulate and transfer the output text attributes (such as formal v.s. informal style, positive v.s. negative sentiment), (3) we should remove social bias and abusive content. Furthermore, I would also like to incorporate personal information such as gender, age, social context, and background knowledge to generate text suitable to individual users.

## References

- [1] **Rui Zhang**, Honglak Lee, and Dragomir R. Radev. Dependency sensitive convolutional neural networks for modeling sentences and documents. In *NAACL*, 2016.
- [2] **Rui Zhang**, Cícero Nogueira dos Santos, Michihiro Yasunaga, Bing Xiang, and Dragomir Radev. Neural coreference resolution with deep biaffine attention by joint mention detection and mention clustering. In *ACL*, 2018.
- [3] **Rui Zhang**, Honglak Lee, Lazaros Polymenakos, and Dragomir Radev. Addressee and response selection in multi-party conversations with speaker interaction rnns. In *AAAI*, 2018.
- [4] Catherine Finegan-Dollak, Reed Coke, **Rui Zhang**, Xiangyi Ye, and Dragomir Radev. Effects of creativity and cluster tightness on short text clustering performance. In *ACL*, 2016.
- [5] Michihiro Yasunaga, **Rui Zhang**, Kshitij Meelu, Ayush Pareek, Krishnan Srinivasan, and Dragomir Radev. Graph-based neural multi-document summarization. In *CoNLL*, 2017.
- [6] Michihiro Yasunaga, Jungo Kasai, **Rui Zhang**, Alexander R Fabbri, Irene Li, Dan Friedman, and Dragomir R Radev. Scisummnet: A large annotated corpus and content-impact models for scientific paper summarization with citation networks. In *AAAI*, 2019.
- [7] **Rui Zhang** and Joel Tetreault. This email could save your life: Introducing the task of email subject line generation. In *ACL*, 2019.
- [8] **Rui Zhang**, Tao Yu, He Yang Er, Sungrok Shim, Eric Xue, Xi Victoria Lin, Tianze Shi, Caiming Xiong, Richard Socher, and Dragomir Radev. Editing-based sql query generation for cross-domain context-dependent questions. In *EMNLP*, 2019.
- [9] Tao Yu, **Rui Zhang**, Michihiro Yasunaga, Yi Chern Tan, Xi Victoria Lin, Suyi Li, Heyang Er, Irene Li, Bo Pang, Tao Chen, Emily Ji, Shreya Dixit, David Proctor, Sungrok Shim, Jonathan Kraft, Vincent Zhang, Caiming Xiong, Richard Socher, and Dragomir Radev. SPaC: Cross-domain semantic parsing in context. In *ACL*, 2019.
- [10] Tao Yu, **Rui Zhang**, He Yang Er, Suyi Li, Eric Xue, Bo Pang, Xi Victoria Lin, Yi Chern Tan, Tianze Shi, Zihan Li, Youxuan Jiang, Michihiro Yasunaga, Sungrok Shim, Tao Chen, Alexander Fabbri, Zifan Li, Luyao Chen, Yuwen Zhang, Shreya Dixit, Vincent Zhang, Caiming Xiong, Richard Socher, Walter Lasecki, and Dragomir Radev. CoSQL: A conversational text-to-sql challenge towards cross-domain natural language interfaces to databases. In *EMNLP*, 2019.
- [11] Tao Yu, **Rui Zhang**, Kai Yang, Michihiro Yasunaga, Dongxu Wang, Zifan Li, James Ma, Irene Li, Qingning Yao, Shanelle Roman, Zilin Zhang, and Dragomir Radev. Spider: A large-scale human-labeled dataset for complex and cross-domain semantic parsing and text-to-SQL task. In *EMNLP*, 2018.
- [12] Catherine Finegan-Dollak, Jonathan K. Kummerfeld, Li Zhang, Karthik Ramanathan, Sesh Sadasivam, **Rui Zhang**, and Dragomir Radev. Improving text-to-SQL evaluation methodology. In *ACL*, 2018.
- [13] Tao Yu, Zifan Li, Zilin Zhang, **Rui Zhang**, and Dragomir Radev. TypeSQL: Knowledge-based type-aware neural text-to-SQL generation. In *NAACL*, 2018.
- [14] Tao Yu, Michihiro Yasunaga, Kai Yang, **Rui Zhang**, Dongxu Wang, Zifan Li, and Dragomir Radev. SyntaxSQLNet: Syntax tree networks for complex and cross-domain text-to-SQL task. In *EMNLP*, 2018.
- [15] **Rui Zhang**, Caitlin Westerfield, Sungrok Shim, Garrett Bingham, Alexander Fabbri, William Hu, Neha Verma, and Dragomir Radev. Improving low-resource cross-lingual document retrieval by reranking with deep bilingual representations. In *ACL*, 2019.
- [16] Douglas W. Oard, Petra Galuščáková, Kathleen McKeown, Marine Carpuat, Mohamed Elbadrashiny, Ramy Eskander, Kenneth Heafield, Efsun Kayi, Chris Kedzie, Smaranda Muresan, Suraj Nair, Xing Niu, Dragomir Radev, Anton Ragni, Han-Chin Shing, Yan Virin, Weijia Xu, **Rui Zhang**, Elena Zotkina, Joseph Barrow, and Mark Gales. Surprise languages: Rapid-response cross-language ir. In *EVIA*, 2019.