# Research Statement

## Sida I. Wang

My research has focused on how to achieve a level of data efficiency in machine learning methods to allow their application in interactive, low-data contexts such as conversational user interfaces (CUIs). While there have been many great recent advances in structured machine learning — the setting where the output space is rich, structured objects — and these advances have allowed us to better tackle important practical problems such as speech recognition, machine translation and parsing, progress has considerably come from demanding ever large supervised training data sets. Researchers have failed to sufficiently address this issue of data efficiency, one that is especially important in CUIs, where semantically annotated data is sparse and adaptation to diverse users and new use cases has to be rapid.
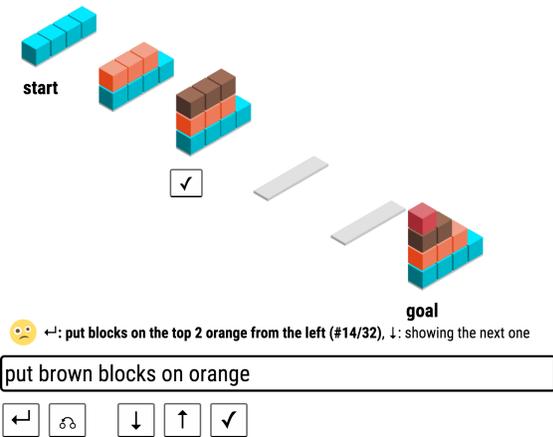
One way of improving data efficiency is through *interaction* — an important ingredient in human language learning — and I use interaction to improve machine language learning (section 1). Furthermore, I also seek to better understand general ML methods through both theoretical and experimental analysis. Understanding data efficiency is a major theme in my research, and often a better understanding of this issue lead to progress in practical problems (section 2).
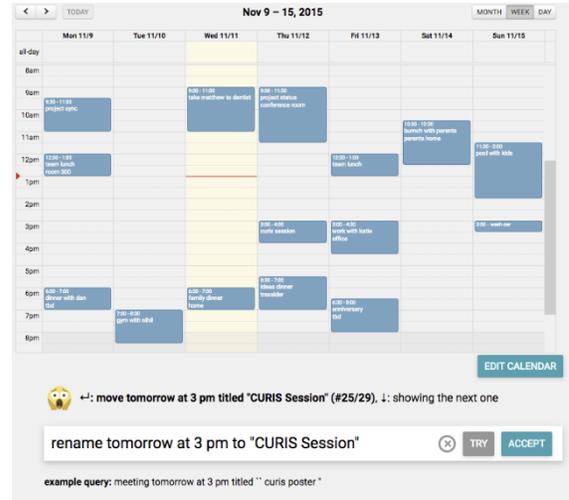
## 1  Interactive language learning

As mobile usage has grown to exceed PC usage and speech recognition is on track to become more usable than keyboards, building good CUIs is clearly the next step. However, high-profile commercial efforts such as Siri and Alexa fall far short of the potential of CUIs. Even with millions of users and lots of data on how people use them, these systems do not get noticeably better over time. My research contributes to the effort towards better CUIs by leveraging user interactions to build more capable and more usable systems that are adaptive.

Specifically to language learning, interaction is an important ingredient for both first and second language acquisition. For first language acquisition, children have diminished ability to acquire language by passively watching TV. For second language acquisition, a popular hypothesis states that *"we acquire by understanding language that contains structure a bit beyond our current level of competence"* and that we obtain this understanding from *"context or extra-linguistic information"*. In contrast, the popular formulation of semantic parsing in NLP provides natural language utterances paired with formal semantic representations without any extra information. In this case, a lot of inference is required from the ML algorithms. Instead of hoping that generic ML algorithms can soon perform super-human inferences for language learning, we can level the playing field by giving our machine language learner an interactive linguistic environment, much like humans have, instead of just a static dataset. In an interactive environment, feedback can be tailored to the current system state.

**1.1  Learning language games**  I studied a setting where language learning starts from scratch with few assumptions, where the system continuously learn from its users through interaction and where rapid learning is key to success — this work was recognized as an outstanding paper at ACL [WLM16]. An inspiration for this setting is Wittgenstein who famously argued that *language derives its meaning from use* and introduced the concept of a *language game* to illustrate the fluidity and purpose-orientedness of language. More concretely, we operationalize and explore the

**(a)** SHRDLURN

**(b)** SCHEDULURN

**Figure 1:** 1a: A pilot for learning language through user interaction. The system attempts an action in response to a user instruction and the user indicates whether it has chosen correctly. This feedback allows the system to learn word meaning and grammar. 1b: the interface for interactive learning in the calendars domain.

idea of language games in a learning setting. In this setting, the two parties do not initially speak a common language, but nonetheless need to collaboratively accomplish a goal. Specifically, we created a game called SHRDLURN, in homage to the seminal work of Winograd (1972). As shown in Figure 1a, the objective is to transform a start state into a goal state, but the only action the human can take is entering an utterance. The computer parses the utterance and produces a ranked list of possible interpretations according to its current model. The human scrolls through the list and chooses the intended one, simultaneously advancing the state of the blocks and providing feedback to the computer. Both the human and the computer wish to reach the goal state (only known to the human) with as little scrolling as possible.

**1.2 The setting** We model the computer as a semantic parser, which maps natural language utterances (e.g., *'remove red'*) into logical forms (e.g., `remove(with(red))`). The semantic parser has no seed lexicon and no annotated logical forms, so it initially generates many candidate logical forms. Based on the human's feedback, it performs online gradient updates on the parameters corresponding to simple and generic lexical features.

It is crucial that the computer learns quickly, or the users are frustrated and the system is less usable. Ultimately we show that good learning speed can be achieved, and can be improved further by modeling *pragmatics*, where the computer also accounts for what the human thinks the computer understands in a probabilistic way.

What is special is the real-time nature of learning, in which the human also learns and adapts to the computer. While the human can teach the computer any language—in our pilot, Mechanical Turk users tried English, Arabic, Polish, and a custom programming language—a good human player will choose to use utterances that the computer is more likely to learn quickly. Thus interaction makes the learning algorithm more data efficient, where it can get away with doing less inference.

**1.3  Current and future directions**   Natural language can often express a very large action space succinctly, often with ambiguities that require context to resolve. While the executable formal language can also be very expressive, the correspondence between the formal language and the natural language tends to become looser as the action space gets more complex. Bridging this gap require richer interactive supervision and I am working on using structured supervision signals such as definitions and demonstrations.

## 2   Machine learning research

**Regularizing with generative models**   It is well-known that generative models are more data efficient, but discriminative models are more flexible and less sensitive to misspecification. I proposed a model, called NBSVM, which uses a generative method to regularize a discriminative method. NBSVM achieved better performance on sentiment analysis and document classification than neural models while being orders of magnitude faster. This method improved the data efficiency of discriminative models and it has been widely cited and implemented [WM12].

**Dropout**   Dropout is a simple and effective method for training neural networks to achieve better data efficiency. It has been very popular because it seems to work in many different models. While the experimental evidence is strong, and while the original work provided some intuitive explanations, a more rigorous analysis was missing. I provided an explanation by approximating the input to each unit as a Gaussian random variable and empirically verified the assumptions [WM13]. Extending that, we derived an explicit regularizer using the delta method, and showed that dropout can be seen as adaptive regularization [WWL13]. Experimentally, we showed that the resulting dropout regularizer improves the performance on structured prediction problems such as Named Entity Recognition [WWW$^+$13].

**MAP inference**   This line of work considers maximum a posteriori (MAP) inference in Markov random fields [WFLM14, FWLM14, FW14]. We studied low-rank relaxations that interpolate between the discrete problem and its full-rank semidefinite relaxation and developed new theoretical bounds on the effect of rank of the relaxation. In practice, we show two algorithms for optimizing the low-rank objectives which are simple to implement, enjoy ties to the underlying theory, and outperform existing approaches on benchmark MAP inference tasks.

**Method of Moments**   Like maximum likelihood, method of moments can be used to estimate parameters of statistical models. While it is usually straightforward to apply the expectation-maximization method – based on maximum likelihood, applying method of moments to various kinds of latent variables models is non-trivial even for experts. I proposed and analyzed a framework for applying method of moments to any mixture model whose moments can be expressed as a polynomial [WCL15].

I want to develop a general framework for interactive language learning, that can easily be used across different domains (figure 1b). My focus is on developing new types of supervision, executable meaning representations and learning algorithms that lead to more data efficiency. I am also interested in formalizing and studying how much data efficiency can we gain by being interactive or by modelling pragmatics.

Looking forward, I believe CUIs have to learn through interaction with users, and become better over time. CUIs have the potential to replace GUIs and scripting for many tasks, and doing so can bridge the great digital divide of skills and enable all of us to better make use of computers.

# References

[FW14] R. Frostig and S. I. Wang. A sub-constant improvement in approximating the positive semidefinite Grothendieck problem. *arXiv preprint arXiv:1408.2270*, 2014.

[FWLM14] R. Frostig, S. I. Wang, P. Liang, and C. Manning. Simple MAP inference via low-rank relaxations. In *Advances in Neural Information Processing Systems (NIPS)*, 2014.

[WCL15] S. I. Wang, A. Chaganty, and P. Liang. Estimating mixture models via mixture of polynomials. In *Advances in Neural Information Processing Systems (NIPS)*, 2015.

[WFLM14] S. I. Wang, R. Frostig, P. Liang, and C. Manning. Relaxations for inference in restricted Boltzmann machines. In *International Conference on Learning Representations Workshop (ICLR)*, 2014.

[WLM16] S. I. Wang, P. Liang, and C. Manning. Learning language games through interaction. In *Association for Computational Linguistics (ACL)*, 2016.

[WM12] S. I. Wang and C. Manning. Baselines and bigrams: Simple, good sentiment and text classification. In *Association for Computational Linguistics (ACL)*, 2012.

[WM13] S. I. Wang and C. Manning. Fast dropout training. In *International Conference on Machine Learning (ICML)*, pages 118–126, 2013.

[WWL13] S. Wager, S. I. Wang, and P. Liang. Dropout training as adaptive regularization. In *Advances in Neural Information Processing Systems (NIPS)*, 2013.

[WWW+13] S. I. Wang, M. Wang, S. Wager, P. Liang, and C. Manning. Feature noising for log-linear structured prediction. In *Empirical Methods in Natural Language Processing (EMNLP)*, 2013.