

---

# Listen, Interact and Talk: Learning to Speak via Interaction

---

Haichao Zhang, Haonan Yu, and Wei Xu  
Baidu Research - Institute of Deep Learning  
Sunnyvale, CA 94089  
{zhanghaichao, haonanyu, xuwei06}@baidu.com

## Abstract

One of the long-term goals of artificial intelligence is to build an agent that can communicate intelligently with human in natural language. Most existing work on natural language learning relies heavily on training over a pre-collected dataset with annotated labels, leading to an agent that essentially captures the statistics of the fixed external training data. As the training data is essentially a static snapshot representation of the knowledge from the annotator, the agent trained this way is limited in adaptiveness and generalization of its behavior. Moreover, this is very different from the language learning process of humans, where language is acquired during communication by taking speaking action and learning from the consequences of speaking action in an interactive manner. This paper presents an interactive setting for grounded natural language learning, where an agent learns natural language by interacting with a teacher and learning from feedback, thus learning and improving language skills while taking part in the conversation. To achieve this goal, we propose a model which incorporates both imitation and reinforcement by leveraging jointly sentence and reward feedbacks from the teacher. Experiments are conducted to validate the effectiveness of the proposed approach.

## 1 Introduction

Natural language is the one of the most natural form of communication for human, and therefore it is of great value for an intelligent agent to be able to leverage natural language as the channel to communicate with human as well. Recent progress on natural language learning mainly relies on supervised training with large scale training data, which typically requires a huge amount of human labor for annotating. While promising performance has been achieved in many specific applications regardless of the labeling effort, this is very different from how humans learn. Humans act upon the world and learn from the consequences of their actions [Skinner, 1957]. For mechanical actions such as movement, the consequences mainly follow geometrical and mechanical principles, while for language, humans act by speaking and the consequence is typically response in the form of verbal and other behavioral feedbacks (*e.g.*, nodding) from conversation partners. These feedbacks typically contain informative signal on how to improve the language skills in subsequent conversions and play an important role in human’s language acquisition process [Petursdottir and Mellor, 2016, Kuhl, 2004, Weston, 2016].

The language acquisition process of a baby is both impressive as a manifestation of human intelligence and inspiring for designing novel settings and algorithms for computational language learning. For example, baby interacts with people and learn through mimicking and feedbacks [Kuhl, 2004, Skinner, 1957]. For learning to speak, baby initially performs verbal action by mimicking his conversational partner (*e.g.* parent) and masters the skill of generating a word (sentence). He could also possibly pick up the association of a word with a visual image when his parents saying “*this is apple*” while pointing to an apple or an image of it. Later, one can ask the baby question like “*what is this*” while pointing

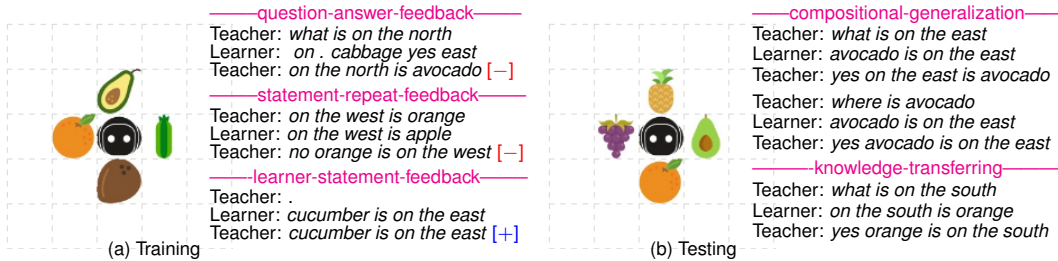


Figure 1: **Interactive language learning example.** (a) During training, teacher interacts in natural language with learner about objects. The interactions are in the form of (1) question-answer-feedback, (2) statement-repeat-feedback, and (3) statement from learner and then feedback from teacher. Certain forms of interactions may be excluded for certain set of object-direction combinations or objects (referred to as *inactive combinations/objects*) during training. For example, the combination of {*avocado, east*} does not appear in question-answer sessions; the object *orange* never appears in question-answer sessions but only in statement-repeat sessions. Teacher provides both sentence feedback as well as reward signal (denoted as  $[+]$  and  $[-]$  in the figure). (b) During testing, teacher can ask question about objects around, including questions involving *inactive combinations/objects* that have never been asked before, e.g., questions about the combination of {*avocado, east*} and questions about *orange*. This testing setup involves *compositional generalization* and *knowledge transferring* settings and is used for evaluating the proposed approach (c.f. Section 4).

to an object, and provides the correct answer if the baby doesn't respond or responds incorrectly, which is typical in the initial stage. One can also provide at the same time a verbal confirmation (e.g. "yes/no") with a nodding/smile/kiss/hug when he answers correctly as a form of encouragement feedback. From a baby's perspective, the way to learn the language is by making verbal utterances to parent and adjusting his verbal behavior according to the corrections/confirmation/encouragement from parent.

This example illustrates that the language learning process is inherently *interactive*, a property which is potentially difficult to be captured by a static dataset as used in the conventional supervised learning setting. Inspired by baby's language learning process, we present a novel interactive setting for grounded natural language learning, where the teacher and the learner can interact with each other in natural languages as shown in Figure 1. In this setting, there is no direct supervisions to guide the behavior of the learner as in the supervised learning setting. Instead, the learner has to *act in order to learn*, i.e., engaging in the conversation with currently acquired speaking skills to obtain feedbacks from the dialogue partner, which provide learning signals for further improvement on the conversation skills.

To leverage the feedbacks for learning, it is tempting to mimic the teacher directly (e.g., using a language model). While this is a viable approach for learning how to speak, the agent trained by pure imitation is not necessarily able to converse adaptively within context due to the negligence of the reinforcement signal. An example is that it is hard to make a successful conversation with a well-trained parrot, which is only good at mimicking. The reason is that the learner is mimicking from a third person perspective [Stadie et al., 2017], mimicking the teacher who is conversing with it, thus certain words in the sentences from the teacher such as "yes/no" and "you/I" might need to be removed/adapted due to the change of perspective from teacher to learner. This cannot be achieved with imitation only. On the other hand, it is also challenging to generate appropriate conversational actions using purely the reinforcement signal without imitation. The fundamental reason is the inability of speaking, thus the probability of generating a sensible sentence by randomly uttering is low, let alone that of a proper one. This is exemplified by the fact that babies don't fully develop their language capabilities without the ability to hear, which is one of the most important channels for language-related imitation.

In this paper, we propose a *joint imitation and reinforcement* approach for interactive language learning. The proposed approach leverages both verbal and encouragement feedbacks from the teacher for joint learning, thus overcoming the difficulties encountered with either only imitation or reinforcement. The contributions of this paper can be therefore summarized as the following:

- We present a novel human-like interaction-based grounded language learning setting where language is learned by interacting with the environment (teacher) in natural language.
- We present a grounded natural language learning approach under the interactive setting by leveraging feedbacks from the teacher during interaction through joint imitation and reinforcement.

To the best of our knowledge, this is the first work on using imitation and reinforcement jointly for grounded natural language learning in an interactive setting.

The remainder of the paper is structured as follows. In Section 2, we make a brief review of related work on natural language learning. Section 3 introduces the formulation of the interaction-based natural language learning problem, followed with detailed explanation of the proposed approach. Experiments are carried out in Section 4 to show the language learning ability of the proposed approach in the interactive setting. Finally, we conclude the paper in Section 5.

## 2 Related Work

Deep network based language learning has received great success recently and has been applied in different applications, for example, machine translation [Sutskever et al., 2014], image captioning/visual question answering [Mao et al., 2015, Vinyals et al., 2015, Antol et al., 2015] and dialogue response generation [Vinyals and Le, 2015, Wen et al., 2015]. For training, a large amount of training data containing source-target pairs is needed, typically requiring a significant amount of efforts to collect. This setting essentially captures the statistics of the training data and does not respect the interactive nature of language learning thus is very different from how humans learn.

While conventional language model is trained in a supervised way, there are some recent works using reinforcement learning for training. These works mainly target at the problem of tuning the performance of a language model pre-trained in a supervised way according to a specific reward function which is either directly the evaluation metric such as standard BLEU core [Ranzato et al., 2016, Bahdanau et al., 2017], manually designed function [Li et al., 2016] or metric learned in an adversarial setting [Yu et al., 2017, Li et al., 2017b], which is non-differentiable, leading to the usage of reinforcement learning. Different from them, our main focus is on the possibility of language learning in an interactive setting and required model designs, rather than optimizing a particular model output towards a specific evaluation metric.

There are some recent works on learning to communicate [Foerster et al., 2016, Sukhbaatar et al., 2016] and the emergence of language [Lazaridou et al., 2017, Mordatch and Abbeel, 2017]. The emerged language need to be interpreted via post-processing [Mordatch and Abbeel, 2017]. Differently, we aim to achieve natural language learning from both perspectives of understanding and generation (*i.e.*, *speaking*), thus the speaking action of the agent is readily understandable without any post-processing. There are also works on dialogue learning using a guesser/responser setting where the guesser tries to achieve the final goal (*e.g.*, classification/localization) by collecting additional information through asking questions to the responser [Strub et al., 2017, Das et al., 2017]. These works try to optimize the question to be asked in order to help the guesser to achieve the final guessing goal. Thus the focus is very different from our goal of language learning through interactions with teacher.

Our work is also related to reinforcement learning based control with natural language action space [He et al., 2016] in the sense that our model also outputs action in natural language space. We also shares similar motivation with [Weston, 2016, Li et al., 2017a], where language learning through textual dialogue has been explored. However, in these works [He et al., 2016, Weston, 2016, Li et al., 2017a] a set of candidate sequences is provided and the action required is selecting one from the candidate set, thus is essentially a *discrete control* problem. In contrast, our model achieves sentence generation through control in a *continuous space*, with a potentially infinite sized action space consisting of all possible sequences.

## 3 Interaction-based Language Learning

We will introduce the proposed interaction-based natural language learning approach in this section. The goal is to design a learning agent<sup>1</sup> that can learn to converse by interacting with the teacher, which can be either a virtual teacher or a human (*c.f.* Figure 1~2). At time step  $t$ , according to a visual image  $\mathbf{v}$ , teacher generates a sentence  $\mathbf{w}^t$  which can be a question (*e.g.*, “*what is on the east*”, “*where is apple*”), a statement (*e.g.*, “*banana is on the north*”), or an empty sentence (denoted as “.”). The learner takes teacher’s sentence  $\mathbf{w}^t$  and the visual content  $\mathbf{v}$ , and produces a sentence response  $\mathbf{a}^t$  to the teacher. The teacher will then provide feedbacks to the learner according to its response in the form of both sentence  $\mathbf{w}^{t+1}$  and reward  $r^{t+1}$ . The sentence  $\mathbf{w}^{t+1}$  represents verbal feedback from teacher (*e.g.*, “*yes on the east is cherry*”, “*no apple is on the east*”) and  $r^{t+1}$  models the non-verbal confirmative feedback such as nodding/smile/kiss/hug, which also appears

---

<sup>1</sup>We use the term *agent* interchangeably with *learner* according to context in the paper.

naturally during interaction. The problem is therefore to design a model that can learn grounded natural language from teacher’s sentences and reward feedbacks. While it might look promising to formulate the problem as supervised training by learning from the subset of sentences from teacher with only positive rewards, this approach won’t work because of the difficulties due to the changed of perspective [Stadie et al., 2017] as mentioned earlier. Our formulation of the problem as well as the details of the proposed approach are presented in the sequel.

### 3.1 Problem Formulation

A response from the agent can be modeled as a sample from a probability distribution over the possible output sequences. Specifically, for one episode, given the visual input  $\mathbf{v}$  and textual input  $\mathbf{w}^{1:t}$  from teacher upto time step  $t$ , the response  $\mathbf{a}^t$  from the agent can be generated by sampling from a policy distribution  $p_\theta^R(\cdot)$  of the speaking action:

$$\mathbf{a}^t \sim p_\theta^R(\mathbf{a}|\mathbf{w}^{1:t}, \mathbf{v}). \quad (1)$$

The agent *interacts* with teacher by outputting the utterance  $\mathbf{a}^t$  and receives the *feedbacks* from teacher at time step  $t + 1$  as  $\mathcal{F} = \{\mathbf{w}^{t+1}, r^{t+1}\}$ .  $\mathbf{w}^{t+1}$  is in the form of a sentence which represents a verbal confirmation/correction in accordance with  $\mathbf{w}^t$  and  $\mathbf{a}^t$ , with prefixes (*yes/no*) added with a probability of half (*c.f.* Figure 1~2). Reward  $r^{t+1}$  is a scalar-valued feedback with positive value as encouragement while negative value as discouragement according to the correctness of the agent utterance  $\mathbf{a}^t$ . The task of interaction-based language learning can be stated as *learning by conversing with teacher and improving from teacher’s feedbacks*  $\mathcal{F}$ . Mathematically, we formulate the problem as the minimization of a cost function as follows:

$$\mathcal{L}_\theta = \mathcal{L}_\theta^I + \mathcal{L}_\theta^R = \underbrace{\mathbb{E}_S \left[ - \sum_t \log p_\theta^I(\mathbf{w}^{t+1}|\mathbf{w}^{1:t}, \mathbf{v}) \right]}_{\text{Imitation}} + \underbrace{\mathbb{E}_{p_\theta^R} \left[ - \sum_t [\gamma]^t \cdot r^{t+1} \right]}_{\text{Reinforce}}, \quad (2)$$

where  $\mathbb{E}_S(\cdot)$  is the expectation over all the sentence sequences  $S$  generated from teacher,  $r^{t+1}$  is the immediate reward received at time step  $t + 1$  after taking speaking action following policy  $p_\theta^R(\cdot)$  at time step  $t$  and  $\gamma$  is the reward discounting factor.  $[\gamma]^t$  is used to denote the exponentiation over  $\gamma$  to differentiate it with superscript indexing. As for both components, the training signal is obtained via *interaction* with the teacher, we termed this task as *interaction-based language learning*. For the imitation part, it essentially learns from teacher’s verbal response  $\mathbf{w}^{t+1}$ , which can only be obtained as a consequence of its speaking action. For the reinforce part, it learns from teacher’s reward signal  $r^{t+1}$ , which is also obtained after taking the speaking action and received at the next time step. The proposed interactive language learning formulation integrates two components which can fully leverage the feedbacks appearing naturally during conversational interaction:

- **Imitation** plays the role of learning a grounded language model by observing teacher’s behaviors during conversation with the learner itself. This enables the learner to have the basic ability to speak within context. The training data here are only the sentences from teacher, without any explicit labeling of ground-truth and is a mixture of expected correct response and others. The way of training is by *predicting the future*. More specifically, the model is predicting the next future word at word level and predicting the next sentence at sentence level. Another important point is that it is in effect *third person imitation* [Stadie et al., 2017], as the learner is imitating the teacher who is conversing with it, rather than another expert student who is conversing with teacher.
- **Reinforce**<sup>2</sup> leverages the confirmative feedbacks from the teacher for learning to converse properly by adjusting the action policy distribution. It enables the learner to use the acquired speaking ability and adapt it according to feedbacks. Here we have the learning signal in the form of reward. This is analogous to baby’s language learning process, who uses the acquired language skills by trial and error with parents and improves according to the encouragement feedbacks.

Note that while imitation and reinforce are represented as two separate components in Eqn.(2), they are tied via parameter sharing in order to fully leverage both forms of training signals. This form of joint learning is crucial for achieving successful language learning, compared with approaches with only imitation or reinforce which are less effective, as verified by experiments in Section 4.

<sup>2</sup>Reinforce denotes the module that learns from the reinforcement/encouragement signal throughout the paper and should be differentiated with the REINFORCE algorithm in the literature [Sutton and Barto, 1998].

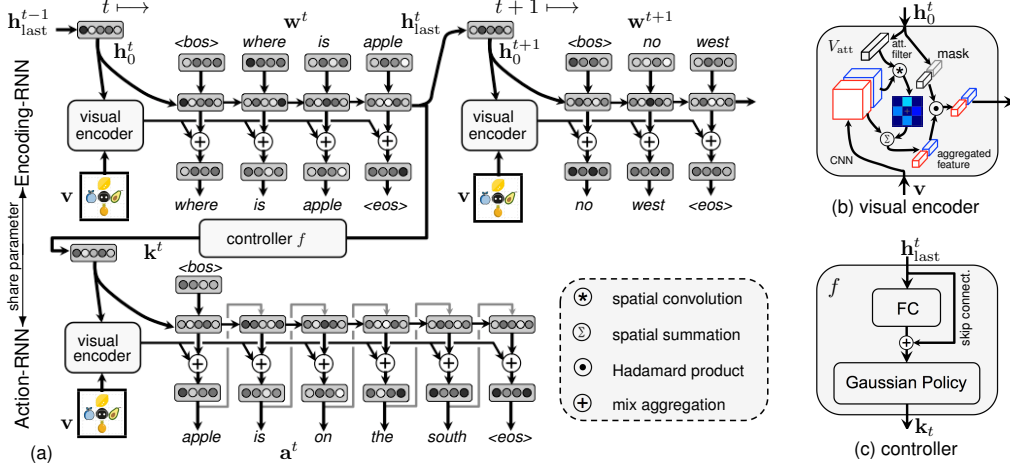


Figure 2: **Network structure.** (a) Illustration of the network structure with sample inputs. (b) Visual encoder network  $V_{att}(\cdot)$ . Visual image is encoded by a CNN and spatially aggregated to a vector with an attention map. The attention map is obtained by convolving the feature map from CNN with a spatial filter generated from the initial state  $h_0^t$ . A binary mask generated from  $h_0^t$  is applied to the spatially aggregated vector to produce the final visual feature vector. At time step  $t$ , the encoding-RNN takes teacher’s sentence (“where is apple”) and the visual feature vector from the visual encoder  $V_{att}(\cdot)$  as inputs. The last state of the encoding-RNN  $h_{last}^t$  is passed through a controller  $f(\cdot)$  to the action-RNN for response generation. Parameters are shared between encoding-RNN and action-RNN. During training, the RNN is trained by predicting next words and next sentences. (c) Controller network with a residue control module followed by a Gaussian Policy module (c.f. Sec. 3.2.2).

### 3.2 Approach

A hierarchical Recurrent Neural Network is used for capturing the sequential structure both across sentences and within a sentence [Yu et al., 2016, Serban et al., 2016], as shown in Figure 2(a). At time-step  $t$ , an encoding-RNN encodes the input sentence  $w^t$  from teacher as well as history information into a state vector  $h_{last}^t$ , which is passed through an action controller  $f(\cdot)$  to produce a control vector  $k^t$  as input to the action-RNN for generating the response  $a^t$  to the teacher’s sentence. Teacher will generate feedback  $\mathcal{F} = \{w^{t+1}, r^{t+1}\}$  according to both  $w^t$  and  $a^t$ . In addition to being used as input to action controller, the state vector is also passed to the next time step and used as the initial state of the encoding-RNN in the next step (*i.e.*,  $h_0^{t+1} \triangleq h_{last}^t$ ) for learning from  $w^{t+1}$ , thus forming another level of recurrence at the scale of time steps.

#### 3.2.1 Imitation with Hierarchical-RNN-based Language Modeling

The teacher’s way of speaking provides a source for the learner to mimic. One way to learn from this source of information is by predictive imitation. Specifically, for a particular episode, we can represent the probability of the next sentence  $w^{t+1}$  conditioned on the previous sentences  $w^{1:t}$  and current image  $v$  as

$$p_\theta^I(w^{t+1}|w^{1:t}, v) = p_\theta^I(w^{t+1}|h_{last}^t, v) = \prod_i p_\theta^I(w_i^{t+1}|w_{1:i-1}^{t+1}, h_{last}^t, v), \quad (3)$$

where  $h_{last}^t$  is the last state of RNN at time step  $t$  as the summarization of  $w^{1:t}$  (*c.f.* Figure 2), and  $i$  indexes words within a sentence. It is natural to model the probability of the  $i$ -th word in the  $t+1$ -th sentence with an RNN as well, where the sentences up to  $t$  and words up to  $i$  within the  $t+1$ -th sentence we conditioned upon is captured by a fixed-length hidden state vector as  $h_i^{t+1} = \text{RNN}(h_{i-1}^{t+1}, w_i^{t+1})$ , thus

$$p_\theta^I(w_i^{t+1}|w_{1:i-1}^{t+1}, h_{last}^t, v) = \text{softmax}(\mathbf{W}_h h_i^{t+1} + \mathbf{W}_v V_{att}(v, h_0^{t+1}) + \mathbf{b}), \quad (4)$$

where  $\mathbf{W}_h$ ,  $\mathbf{W}_v$  and  $\mathbf{b}$  denote the transformation weight and bias parameters respectively.  $V_{att}(\cdot)$  denotes the visual encoding network with spatial attention incorporated as shown in Figure 2(b).  $V_{att}(\cdot)$  takes the initial RNN state  $h_0^t$  and visual image  $v$  as input. The visual image is first encoded by a CNN to obtain the visual feature map (red cube in Figure 2(b)). The visual feature map is appended with another set of maps with learnable parameters for encoding the directional information (blue cube in Figure 2(b)). This set of feature maps is spatially aggregated to a vector with an attention

map, which is obtained by convolving the feature map with a spatial filter generated from the initial RNN state. An attention mask for emphasizing visual or directional features is generated from  $\mathbf{h}_0^t$  and is applied to the spatially aggregated vector to produce the final feature vector. The initial state of the encoding-RNN is the last state of the previous RNN, *i.e.*,  $\mathbf{h}_0^{t+1} = \mathbf{h}_{\text{last}}^t$  and  $\mathbf{h}_0^0 = \mathbf{0}$ .

The language model trained this way will have the basic ability of producing a sentence conditioned on the input. Therefore, when connecting an encoding-RNN with action-RNN directly, *i.e.*, inputting the last state vector from encoding-RNN into action-RNN as the initial state, the learner will have the ability to generate a sentence by mimicking the way teacher speaks, due to parameter sharing. However, this basic ability of speaking is not enough for the learner to converse properly with teacher, which requires the incorporation of reinforcement signals as detailed in the following section.

### 3.2.2 Learning via Reinforcement for Sequence Actions

An agent generates an action according to  $p_{\theta}^R(\mathbf{a}|\mathbf{w}^{1:t}, \mathbf{v})$ . As sentences  $\mathbf{w}^{1:t}$  can be summarized as the last RNN state  $\mathbf{h}_{\text{last}}^t$ , the action policy distribution can be represented as  $p_{\theta}^R(\mathbf{a}|\mathbf{h}_{\text{last}}^t, \mathbf{v})$ . To leverage the language skill that is simultaneously learned from imitation, we can generate the sentence using a language model shared with imitation, but with a modulated conditional signal via a controller network  $f(\cdot)$  as follows (*c.f.* Figure 2(a, c))

$$p_{\theta}^R(\mathbf{a}^t|\mathbf{h}_{\text{last}}^t, \mathbf{v}) = p_{\theta}^I(\mathbf{a}^t \triangleq \mathbf{w}^{t+1}|f(\mathbf{h}_{\text{last}}^t), \mathbf{v}). \quad (5)$$

The reason for incorporating a controller  $f(\cdot)$  for modulation is that the basic language model only offers the learner the ability to generate a sentence, but not necessarily the ability to respond correctly, or to answer a question from teacher properly. Without any additional module, the agent’s behaviors would be the same as those from teacher because of parameter sharing, thus agent cannot learn to speak correctly in an adaptive manner by leveraging the feedbacks from teacher.

Controller  $f(\cdot)$  is a composite function with two components: (1) a residue structured network for transforming the encoding vector  $\mathbf{h}_{\text{last}}^t$  in order to modify the behavior; (2) a Gaussian policy module for generating a control vector from a Gaussian distribution conditioned on the transformed encoding vector from the residue control network as a form of exploration. In practice, we also incorporate a gradient-stopping layer between the controller and its input, to encapsulate all the modulation ability within the controller.

**Residue Control.** The action controller should have the property that it can pass the input vector to the next module unmodified when desirable while can modify the content of the input vector otherwise. Therefore, a residue structured network is one design satisfying this requirement, with a content modifying vector added to the original input vector (*i.e.*, skip connection) as follows

$$\mathbf{c} = \tau(\mathbf{h}) + \mathbf{h}, \quad (6)$$

where  $\tau(\cdot)$  is a content transformation net and  $\mathbf{c}$  is the generated control vector. The reason for including a skip connection is that it offers the ability to leverage the language model simultaneously learned via imitation for generating sensible sentences and the transformation net  $\tau(\cdot)$  includes learnable parameters for adjusting the behaviors via interactions with the environment and feedbacks from teacher. We implement  $\tau(\cdot)$  as two fully-connected layers with ReLU activation.

**Gaussian Policy.** Gaussian policy net models the output vector as a Gaussian distribution conditioned on the input vector. It takes the generated control vector  $\mathbf{c}$  as input and produces a vector  $\mathbf{k}$  which is used as the initial state of the action-RNN. The Gaussian policy is modeled as follows:

$$p_{\theta}^R(\mathbf{k}|\mathbf{c}) = \mathcal{N}(\mathbf{c}, \mathbf{\Gamma}^T \mathbf{\Gamma}), \quad \mathbf{\Gamma} = \text{diag}[\gamma(\mathbf{c})]. \quad (7)$$

The incorporation of Gaussian policy introduces stochastic unit into the network, thus back-propagation cannot be applied directly. We therefore use policy gradient algorithm for optimization [Sutton and Barto, 1998]. where  $\gamma(\cdot)$  is a sub-network for estimating the standard derivation vector and is implemented using a fully-connected layer with ReLU activation.<sup>3</sup> The vector  $\mathbf{k}$  generated from the controller is then used as the initial state of action-RNN and the sentence output is generated using beam search (*c.f.* Figure 2(a)). For the reward  $r^{t+1}$  in Eqn.(2), we introduce a baseline for reducing variance as  $r^{t+1} - V_v(\mathbf{v})$ , where  $V_v(\cdot)$  represents the value network with parameter vector  $v$  and is estimated by adding to  $\mathcal{L}^R$  an additional value network cost  $\mathcal{L}^V$  as follows

$$\mathcal{L}^V = \mathbb{E}_{p_{\theta}^R}(r^{t+1} + \lambda V_{v^-}(\mathbf{v}^{t+1}) - V_v(\mathbf{v}^t))^2, \quad (8)$$

where  $v$  denotes the set of parameters in the value network and  $V_{v^-}(\cdot)$  denotes the target version of the value network, whose parameter vector  $v^-$  is periodically copied from the training version [Mnih et al., 2013].

<sup>3</sup>In practice, we add a small value (0.01) to  $\gamma(\mathbf{c})$  as a constrain of the minimum standard deviation.

### 3.3 Training

Training involves optimizing the stochastic policy by using the teacher’s feedback  $\mathcal{F}$  as a training signal, obtaining a set of optimized parameters by considering jointly imitation and reinforcement as shown in Eqn.(2). Stochastic gradient descend is used for training the network. For  $\mathcal{L}^I$  from imitation module, we have its gradient as:

$$\nabla_{\theta} \mathcal{L}_{\theta}^I = -\mathbb{E}_S [\nabla_{\theta} \sum_t \log p_{\theta}^I(\mathbf{w}^{t+1} | \mathbf{w}^{1:t}, \mathbf{v})]. \quad (9)$$

Using policy gradient theorem [Sutton and Barto, 1998], we have the following gradient for the reinforce module:

$$\nabla_{\theta} \mathcal{L}_{\theta}^R = -\mathbb{E}_{p_{\theta}^R} [[\nabla_{\theta} \log p_{\theta}^R(\mathbf{k}^t | \mathbf{c}^t) + \nabla_v V_v(\mathbf{v})] \cdot \delta], \quad (10)$$

where  $\delta$  is the td-error defined as  $\delta = r^{t+1} + \gamma V_v(\mathbf{v}) - V_v(\mathbf{v})$ . The algorithm is implemented with deep learning platform PaddlePaddle<sup>4</sup>. We train the network with Adagrad [Duchi et al., 2011] with a batch size of 16 and a learning rate of  $1 \times 10^{-5}$ . Discount factor  $\gamma = 0.99$ . Experience replay is used in practice [Mnih et al., 2013].

## 4 Experimental Results

We evaluate the performance of the proposed approach under several different settings to demonstrate its ability of interactive language learning. For training efficiency, we construct a simulated environment for language learning as shown in Figure 1. We consider the case with four different objects around the learner in each direction ( $S, N, E, W$ ), which are randomly sampled from a set of objects for each session. Within this environment, a teacher interacts with the agent about objects around in three different forms: (1) asking a question as “*what is on the south*”, “*where is apple*” and the agent answers the question; (2) describing objects around as “*apple is on the east*” and agents repeat the statement; (3) saying nothing (“.”) then agent describes objects around and gets a feedback from teacher. The agent receives a positive reward ( $r=+1$ ) if it behaves correctly (generates a correct answer to a question from teacher or produces a correct statement if teacher says nothing) and a negative reward ( $r=-1$ ) otherwise. Reward is used to represent teacher’s non-verbal feedback such as *nodding* as a form of encouragement. Besides reward feedback, teacher also provides a verbal feedback including the expected answer in the form of “*X is on the east*” or “*on the east is X*” and with prefix (“*yes/no*”) added with a probability of half. The speaking action from the agent is correct if it outputs a sentence matches exactly with the expected answer in one of the above forms. There is a possibility for the learner to generate a new correct sentence that beyond teacher’s knowledge. This is not handled in our current work due to the usage of a scripted teacher.

**Language Learning Evaluation.** We first validate the basic language learning ability of the proposed approach under the interactive language learning setting. In this setting, the teacher first generates a sentence for the learner, then the learner will respond, and the teacher will provide feedback in terms of sentence and reward. We compare the proposed approach with two baseline approaches: (1) **Reinforce** which uses directly reinforcement for learning from teacher’s reward feedback [Sutton and Barto, 1998]; (2) **Imitation** which learns by mimicking teacher’s behavior [Sutskever et al., 2014]. Experimental results are shown in Figure 3. It is interesting to note that learning directly from reward feedback only (**Reinforce**) does not lead to successful language acquisition. This is mainly because of the low possibility of generating a sensible sentence by random exploration, and the even lower possibility of generating the correct sentence, thus the received reward can stay at  $-1$ . On the other hand, the **Imitation** approach performs better than **Reinforce**, due to the *speaking* ability it gained through mimicking. The proposed approach achieves reward higher than both compared approaches, due to the effectiveness of the joint formulation, which can fully leverage the feedback signals appeared naturally during conversation for learning. This indicates the effectiveness of the proposed approach for language learning under the interactive setting. Similar behaviors have been observed during testing. We further visualize some examples as shown in Figure 4 along with the generated attention maps. As can be observed from the results, the proposed approach can successfully generate correct attention maps for both *what* and *where* questions. When teacher says nothing (“.”), the agent can generate a statement describing an object around correctly.

<sup>4</sup><https://github.com/PaddlePaddle/Paddle>

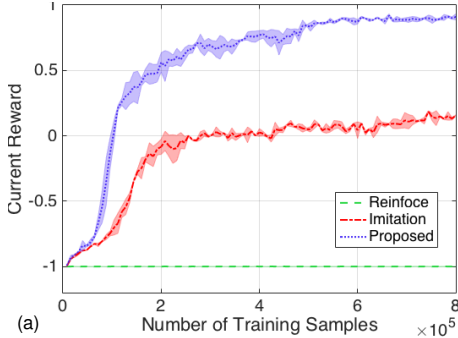


Table 1: Testing Results with Mixed Config.

Settings	Reinforce	Imitation	Proposed
Compositional-gen.	0.0%	83.7%	98.9%
Knowledge-transfer	0.0%	81.6%	97.5%

Table 2: Testing Results with Held-out Config.

Settings	Reinforce	Imitation	Proposed
Compositional-gen.	0.0%	75.1%	98.3%
Knowledge-transfer	0.0%	70.4%	89.0%

Figure 3: **Evaluation results.** (a) Evolution of reward during training. (b) Comparison of the proposed approach with Reinforce and Imitation approaches across different settings and configurations. *Mixed config* denotes the configuration involving interactions with all objects. *Held-out config* denotes the configuration involving interactions with only the objects that are inactive during training.

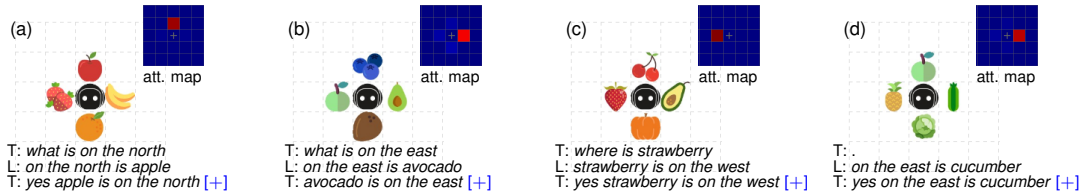


Figure 4: **Example results.** (a-b) *what* questions. (c) *where* question. (d) teacher says nothing (“.”) and the agent is expected to produce a statement. For each example, we show the visual image, the conversion dialogues between teacher and learner, as well as the *attention map* (att. map) generated from the learner when producing the response to teacher (overlaid on top-right). The attention map is rendered as a heat map, with red color indicating large value while blue indicating small value. Grid lines are overlaid on top of the attention map for visualization purpose. The position of the learner is marked with a cross in the attention map (T/L: teacher/learner, [+/-]: positive/negative rewards).

**Zero-shot Dialogue.** An intelligent agent is expected to have an ability to generalize. While this is not the main focus on the paper, we use it as a way to assess the language learning ability of an approach. Experiments are done in following two settings. (1) **Compositional generalization:** the learner interacts with the teacher about objects around during training, but does not have any interaction with certain objects (referred to as *inactive* objects) at particular locations, while in testing the teacher can ask questions about an object regardless of its location. It is expected that a good learner should be able to generalize the concepts it learned about both *objects* and *locations* as well as the acquired conversation skills and can interact successfully in natural language with teacher about novel  $\{object, location\}$  combinations that it never experienced before. (2) **Knowledge transferring:** teacher asks learner questions about the objects around. For certain objects, the teacher only provides descriptions without asking questions during training, while in testing, the teacher can ask questions about any object present in the scene. The learner is expected to be able to transfer the knowledge learned from teacher’s description to generate an answer to teacher’s question about these objects. Experiments are carried out under these two settings for two configurations (*mixed* and *held-out*) and experimental results are summarized in Table 1 and Table 2 respectively. *Mixed configuration* denotes the case in which a mixture of interactions with all objects regardless of whether they are active or inactive during training. *Held-out configuration* denotes the case involving interactions with only the objects that are inactive during training. The results shows that the **Reinforce** approach performs poorly under both settings due to the lack of basic language-related abilities as mentioned in the previous section. The **Imitation** approach performs better than **Reinforce** mainly due to its language speaking ability through mimicking. Note that the held-out configuration is a subset of the mixed-configuration involving only novel objects/combinations, thus is more difficult than the mixed case. It is interesting to note that the proposed approach maintains a consistent behavior under the more difficult held-out configuration and outperforms the other two approaches under both settings, demonstrating its effectiveness in interactive language learning.



## 5 Conclusion

We present an interactive setting for grounded natural language learning and propose an approach that achieves effective interactive natural language learning by fully leveraging the feedbacks that arise naturally during interaction through joint imitation and reinforcement. Experimental results show that the proposed approach provides an effective way for natural language learning in the interactive setting and enjoys desirable generalization and transferring abilities under several different scenarios. As for future work, we would like to explore the direction of explicit modeling of learned knowledge [Yang, 2003] and fast learning about new concepts [Andrychowicz et al., 2016]. Another interesting direction is to connect the language learning task presented in this paper with other heterogeneous tasks such as navigation.

## Acknowledgements

We thank Xiaochen Lian, Zhuoyuan Chen, Yi Yang and Qing Sun for their discussions and comments.

## References

- M. Andrychowicz, M. Denil, S. G. Colmenarejo, M. W. Hoffman, D. Pfau, T. Schaul, and N. de Freitas. Learning to learn by gradient descent by gradient descent. In *NIPS*, 2016.
- S. Antol, A. Agrawal, J. Lu, M. Mitchell, D. Batra, C. L. Zitnick, and D. Parikh. VQA: Visual Question Answering. In *ICCV*, 2015.
- D. Bahdanau, P. Brakel, K. Xu, A. Goyal, R. Lowe, J. Pineau, A. C. Courville, and Y. Bengio. An actor-critic algorithm for sequence prediction. In *ICLR*, 2017.
- A. Das, S. Kottur, J. M. Moura, S. Lee, and D. Batra. Learning cooperative visual dialog agents with deep reinforcement learning. *CoRR*, abs/1703.06585, 2017.
- J. C. Duchi, E. Hazan, and Y. Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research*, 12:2121–2159, 2011.
- J. N. Foerster, Y. M. Assael, N. de Freitas, and S. Whiteson. Learning to communicate with deep multi-agent reinforcement learning. In *NIPS*, 2016.
- J. He, J. Chen, X. He, J. Gao, L. Li, L. Deng, and M. Ostendorf. Deep reinforcement learning with a natural language action space. In *ACL*, 2016.
- P. K. Kuhl. Early language acquisition: cracking the speech code. *Nat Rev Neurosci*, 5(2):831–843, 2004.
- A. Lazaridou, A. Peysakhovich, and M. Baroni. Multi-agent cooperation and the emergence of (natural) language. In *ICLR*, 2017.
- J. Li, W. Monroe, A. Ritter, D. Jurafsky, M. Galley, and J. Gao. Deep reinforcement learning for dialogue generation. In *EMNLP*, 2016.
- J. Li, A. H. Miller, S. Chopra, M. Ranzato, and J. Weston. Learning through dialogue interactions. In *ICLR*, 2017a.
- J. Li, W. Monroe, T. Shi, A. Ritter, and D. Jurafsky. Adversarial learning for neural dialogue generation. *CoRR*, abs/1701.06547, 2017b.
- J. Mao, W. Xu, Y. Yang, J. Wang, Z. Huang, and A. Yuille. Deep captioning with multimodal recurrent neural networks (m-RNN). *ICLR*, 2015.
- V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, and M. Riedmiller. Playing Atari with deep reinforcement learning. In *NIPS Deep Learning Workshop*. 2013.
- I. Mordatch and P. Abbeel. Emergence of grounded compositional language in multi-agent populations. *CoRR*, abs/1703.04908, 2017.
- A. I. Petrusdottir and J. R. Mellor. Reinforcement contingencies in language acquisition. *Policy Insights from the Behavioral and Brain Sciences*, 4(1):25–32, 2016.
- M. Ranzato, S. Chopra, M. Auli, and W. Zaremba. Sequence level training with recurrent neural networks. In *ICLR*, 2016.
- I. V. Serban, A. Sordoni, Y. Bengio, A. C. Courville, and J. Pineau. Building end-to-end dialogue systems using generative hierarchical neural network models. In *AAAI*, 2016.
- B. F. Skinner. *Verbal Behavior*. Copley Publishing Group, 1957.
- B. C. Stadie, P. Abbeel, and I. Sutskever. Third-person imitation learning. In *ICLR*, 2017.
- F. Strub, H. de Vries, J. Mary, B. Piot, A. C. Courville, and O. Pietquin. End-to-end optimization of goal-driven and visually grounded dialogue systems. In *IJCAI*, 2017.
- S. Sukhbaatar, A. Szlam, and R. Fergus. Learning multiagent communication with backpropagation. In *NIPS*, 2016.
- I. Sutskever, O. Vinyals, and Q. V. Le. Sequence to sequence learning with neural networks. In *NIPS*, 2014.
- R. S. Sutton and A. G. Barto. *Reinforcement Learning: An Introduction*. MIT Press, 1998.
- O. Vinyals and Q. V. Le. A neural conversational model. *CoRR*, abs/1506.05869, 2015.

- O. Vinyals, A. Toshev, S. Bengio, and D. Erhan. Show and tell: A neural image caption generator. In *CVPR*, 2015.
- T. Wen, M. Gasic, N. Mrksic, P. Su, D. Vandyke, and S. J. Young. Semantically conditioned LSTM-based natural language generation for spoken dialogue systems. In *EMNLP*, 2015.
- J. Weston. Dialog-based language learning. In *NIPS*, 2016.
- C. D. Yang. *Knowledge and Learning in Natural Language*. Oxford University Press UK, 2003.
- H. Yu, J. Wang, Z. Huang, Y. Yang, and W. Xu. Video paragraph captioning using hierarchical recurrent neural networks. In *CVPR*, 2016.
- L. Yu, W. Zhang, J. Wang, and Y. Yu. SeqGAN: Sequence generative adversarial nets with policy gradient. In *AAAI*, 2017.