

# Placeholder

# Placeholder

## Applying Seq2seq in Building Task-oriented Conversation Systems

### Background

- Logic-based conversation systems in current business solutions can be increasingly complex in structure.
- We aim to build a **simplistic** machine learning model that is able to **learn** from the contents in dialogs.
- We choose **Seq2seq**, which comes with the edge to align up context information in a sequential manner, as the core component.
- Many current applications focus more on building a "chat-bot" to cover generalized topics and mimic human dialogs.
- We choose to build a **task-oriented** conversation system to provide intelligent business specified information to customers, similar to Insurance QA.



Intern: Yiran Luo

Mentor: Lijun Mei, Yipeng Yu

Manager: Shaochun Li



### Challenges

book\_title : andrew\_lang\_prince\_prioig.txt.out  
chapter 1. -lcb- chapter heading picture : pl.jpg -rcb-  
once upon a time there reigned in pantoufria a king and  
with almost everything else to make them happy, they  
this vexed the king more than the queen. However,  
however, she, KAT +++\$+++ It's Shakespeare. Maybe you've heard of  
the king was anxious. BIANCA +++\$+++ Like I'm supposed to know what that  
she did not believe KAT +++\$+++ At least I'm not a clouted fen- sucked  
nothing else. BIANCA +++\$+++ Can't you forget for just one night  
BIANCA +++\$+++ Boney Lowenstein's party is normal,  
KAT +++\$+++ What's the normal?

Source: Facebook Children's Book Test (back) and Comell's Movie-DialogsCorpus (front)

- The essential data for training within such a designated circumstance are highly restricted. Most of existing training conversations are for open domain purposes.

不计免赔是所有的险种都包含了，比如说开车不小心刮蹭了一下，那保险公司在定责任的时候会有一分15%是自己的

就按去年的续

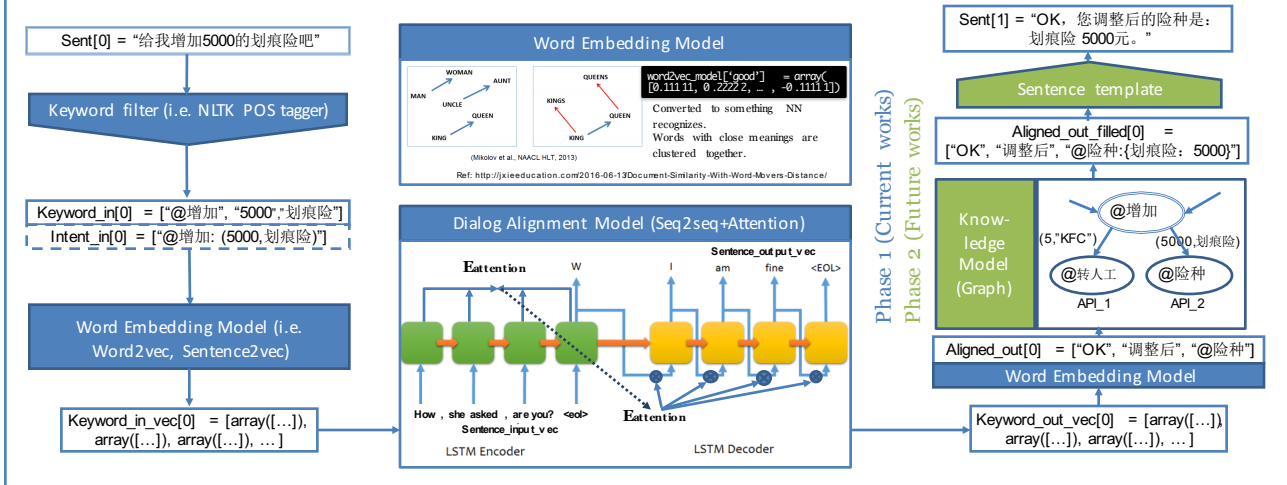
好的，那就按照总金额5200的给您续，对吧？

对

好的，根据北京保险行业协会要求，我们给您发送了一个手机验证码以确认身份信息采集，请您收到后输入手机验证码

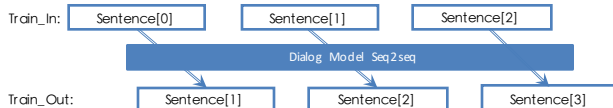
- Open domain data comes with lots of noise. Words such as prepositions are unnecessary in maintaining key information in the context of dialogs.
- Thus we may indefinitely approach our targets, but may never retrieve the exactly same words.

### Framework



### Experiment & Result

- Assumptions:** Since few task-oriented dialog data is available at hand, we simulate such a conversation scenario with open domain data. Test sentences are selected from a single story (as a task frame) and we align up every two neighboring sentences while training the Seq2seq model.



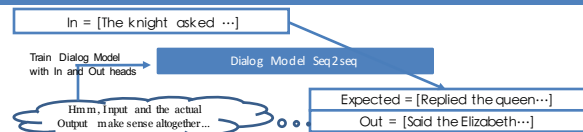
- Experiment settings:** Facebook Children's Book Test Corpus 200k sentences, 2000 sentences as test data, 4-layer NN, word vector dimension 100, hidden layer dimension 200, single GPU, 40+ epochs.
- Word embedding tested with both Word2vec and Sentence2vec. All words are converted to lowercase. Punctuation are limited to " , . . . . ."
- In a single training batch, short sentences are padded with ' ' at the end.
- Implemented in Python with Keras, Theano, gensim and CUDA.

- Result:** Achieved high similarity while aligning short sentences, but not so well at aligning long sentences so far.

output	target	cos-similarity
forget	fiddle-de-dee	0.889915
!	!	0.673258
...	...	...
said	replied	0.761818
the	the	1
prince	queen	0.787364
caddy	;	0.670446
for	for	0.58336
honeycomb	the	0.528529
honeycomb	king	0.419499

- Discussion:** In this experiment, we aim at overfitting instead of avoiding it. This may bring down test accuracy given a similar sentence but with slight differences, or given an exactly existing sentence but from another user.

### Innovations



- FasterConvergence:** An idea somewhat inspired by reinforcement learning. If we find a pair of good-looking aligned sequences (partial or not), we train the pre-trained Seq2seq model with this very pair for multiple epochs, letting the model recognize the good alignment faster.
- For the Word2vec test, quantitatively, we set a good pair to have an average cosine similarity of over 0.8, with over 3 words in a row. Or when the test case is small (<2000 sentences), we can pick the good pairs out by hand.

### Future Works: Knowledge Model

Intent (out/in)	Intent_1	...	@增加	...
Intent_1	{kw_5:0.4, Action2 kw_1:0.1, Action1}	...	0	...
...	...	...	...	...
@险种	0	...	{例:险种:0.83,Add, 人身险:0.17,Delete}	...
@增加	0	...	0	...

- Ideas include using an adjacency matrix to build a frequency-based knowledge model, and introducing intent keywords to carry user-specific information throughout knowledge retrieval/update.