



Molweni: A Challenge Multiparty Dialogues-based Machine Reading Comprehension Dataset with Discourse Structure

Jiaqi Li, Ming Liu, Min-Yen Kan, Zihao Zheng, Zekun Wang,
Wenqiang Lei, Ting Liu, Bing Qin

Harbin Institute of Technology, National University of Singapore, Peng Cheng Laboratory



▼ INTRODUCTION

- We present the *Molweni* dataset, a machine reading comprehension (MRC) dataset with discourse structure built over multiparty dialog. Molweni's source samples from the Ubuntu Chat Corpus, including 10,000 dialogs comprising 88,303 utterances. We annotate 30,066 questions on this corpus, including both answerable and unanswerable questions.
- Molweni also uniquely contributes discourse dependency annotations in a modified Segmented Discourse Representation Theory (SDRT)) style for all of its multiparty dialogs, contributing large-scale (78,245 annotated discourse relations) data to bear on the task of multiparty dialog discourse parsing.
- Our experiments show that Molweni is a challenging dataset for current MRC models: BERT-wwm, a current, strong SQuAD 2.0 performer, achieves only 67.7% F_1 on Molweni's questions, a 20+% significant drop as compared against its SQuAD 2.0 performance.

Dialogue 1		
<i>nbx909</i> : how do i find the address of a usb device ?	U_1	Q1: Why does <i>likwidoxigen</i> do a full restart?
<i>likwidoxigen</i> : try taking it out to dinner and do a little wine and dine and it should tell ya	U_2	A1: it re-loads the modules
<i>likwidoxigen</i> : what sort of device ?	U_3	Q2: What does <i>nbx909</i> want to do?
<i>babo</i> : ca n't i just copy over the os and leave the data files untouched ?	U_4	A2: find the address of a usb device
<i>nbx909</i> : only if you do an upgrade	U_5	Q3: How to restart network?
<i>nuked</i> : should i just restart x after installing	U_6	A3: NA.
<i>likwidoxigen</i> : i 'd do a full restart so that it re-loads the modules	U_7	

(a)

(b)

Figure 1. (Dialog 1) A corpus example from Molweni. There are four speakers in the dialog: *nbx909*, *likwidoxigen*, *babo*, and *nuked*. In total, the speakers make seven utterances: U_1 to U_7 . Our annotators proposed three questions against the provided dialog: Q1–3, where Q1 and Q2 are answerable questions, and Q3 is unanswerable. Due to the properties of informal dialog, the instances in our corpus often have grammatical errors.

▼ OVERVIEW OF *MOLWENI*

	Train	Dev	Test	Total
Number of Dialogs	8,771	883	100	9,754
Number of Utterances	77,374	7,823	845	86,042
Number of Questions	24,682	2,513	2,871	30,066

Table1. Overview of Molweni for MRC.

	Train	Dev	Test	Total
Number of Dialogs	9,000	500	500	10,000
Number of Utterances	79,487	4,386	4,430	88,303
Number of Relations	70,454	3,880	3,911	78,245

Table2. Overview of Molweni for DP.

Metric	Number
Average / Maximum number of speakers per dialog	3.51 / 9
Average / Maximum question length (in tokens)	5.91 / 18
Average / Maximum answer length (in tokens)	4.08 / 19
Average / Maximum dialogue length (in tokens)	104.4 / 208
Average / Maximum dialogue length (in utterances)	8.82 / 14
Vocabulary size	24,615
Answerable questions	25,779
Unanswerable questions	4,287

Table3. Detailed statistics for the Molweni corpus.

▼ ANNOTATION FOR MACHINE READING COMPREHENSION

Dataset	Answer type	Dialogue text	Multiparty dialogue	Unanswerable questions	Discourse structure
RACE (Lai et al., 2017)	multiple-choice	✓	✓	✓	✓
NarrativeQA (Kocisky et al., 2018)	abstractive	✓	✓	✓	✓
CoQA (Choi et al., 2018)	abstractive	✓	✓	✓	✓
SQuAD 2.0 (Rajpurkar et al., 2018)	extractive	✓	✓	✓	✓
QuAC (Choi et al., 2018)	extractive	✓	✓	✓	✓
(Ma et al., 2018)	cloze	✓	✓	✓	✓
DREAM (Sun et al., 2019)	multiple-choice	✓	✓	✓	✓
FriendsQA (Yang and Choi, 2019)	extractive	✓	✓	✓	✓
Molweni (Our)	extractive	✓	✓	✓	✓

Table4. Comparison of Molweni with other MRC datasets on answer type, text type (dialogue or written text), multiparty dialogs or not, unanswerable questions, and discourse structure.

Question	Example		Proportion(%)
How	How to do an upgrade?	How can I use this machine?	9.9
Why	Why is it not mounted?	Why does <i>jimcooanct</i> meet the error?	4.3
Who	Who is chart's service customers?	Who is using ubuntu?	4.7
When	When does <i>rhodry</i> have the error?	When is <i>SuperMiguel</i> back?	1.7
Where	Where did <i>earthen</i> write in?	Where is the device?	5.7
What	What does <i>elnomade</i> choose?	What does <i>noone</i> need?	71.7
Others	Does <i>elnomade</i> choose the print?	Which version does <i>xxiao</i> find?	1.9

Table5. Examples of questions in Molweni.

▼ ANNOTATION FOR DISCOURSE PARSING

Dialogue 2

toma:- but its well worth the wait
woodgrain: i have a decently fast p4 should i still be waiting ?
toma:- have you run updatedb before ?
woodgrain: no never before -- but it worked and now i have all the files i need .
woodgrain: i do n't have a path to the jre -- do i need to add it ?
toma:- a path ? ? you compiling somehting ?
woodgrain: do n't need jdk as witnessed by eclipse irc
woodgrain: no i 'm installing this newer ver from the eclipse site .

U_1
 U_2
 U_3
 U_4
 U_5
 U_6
 U_7
 U_8

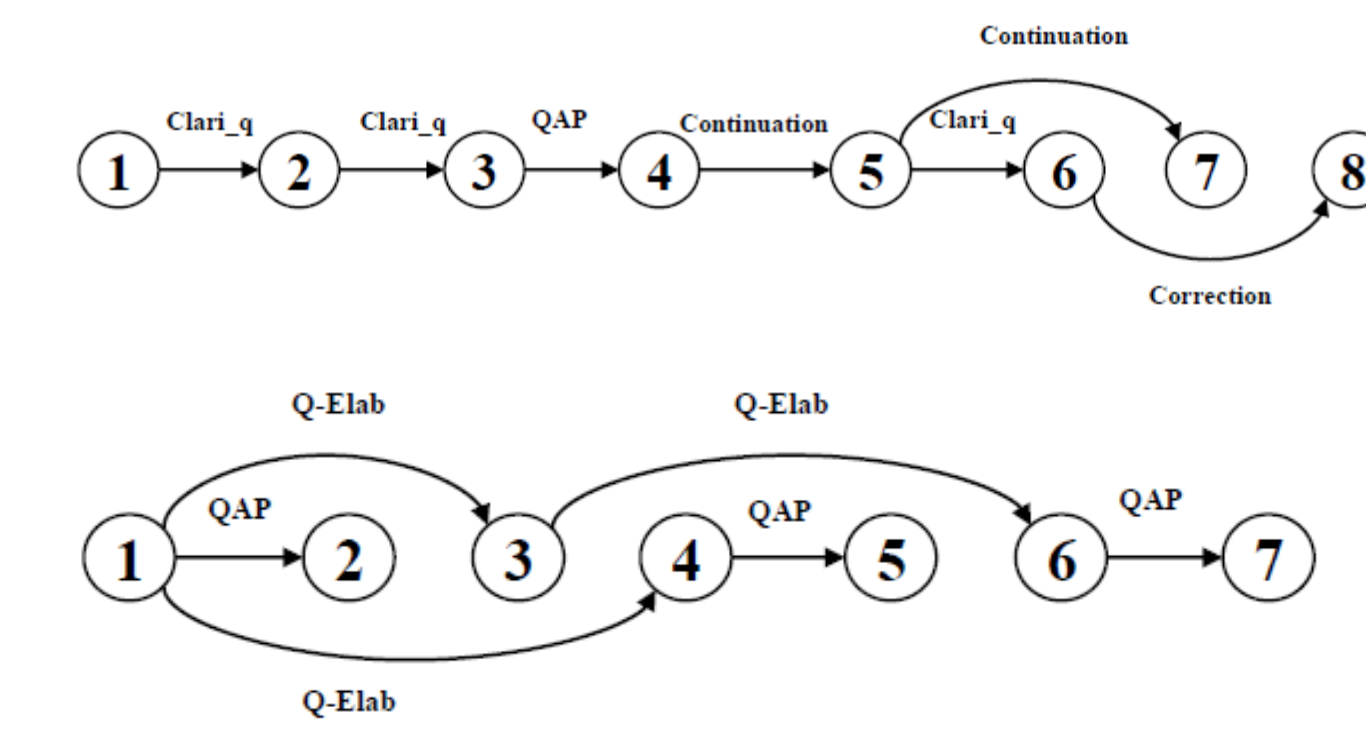


Figure 2. Dialog 2 is a two-party dialog example with eight utterances— U_1 to U_8 —proposed by two speakers: *toma*- and *woodgrain*. The discourse dependency structure and relations for Dialog 2 (Top, two-party) and Dialogue 1 (Bottom, multiparty). Clari_q, QAP, and Q-Elab are respectively short for Clarification question, Question-answer pair, and Question-Elaboration. The label on the link represents the discourse dependency relations between two utterances.

▼ DATA QUALITY

- Manually Check.
- Programmatic Check.
- The FLESS Kappa: 0.91 for discourse dependency links, and 0.56 for both links and discourse relations.

▼ EXPERIMENTAL RESULTS

Method	EM		F1	
	Squad 2.0	Our	Squad 2.0	Our
BERT-base	73.1	45.3	76.2	58.0
BERT-large	80.0	51.8	83.1	65.5
BERT-wwm	86.7	54.7	89.1	67.7
Human performance	86.8	64.3	89.4	80.2
Human-machine gap	0.1	9.6	0.3	12.5

Table6. Results of machine reading comprehension for multiparty dialogs.

Method	Link		Link & Relation	
	STAC	Our	STAC	Our
Deep sequential	73.2	78.1	55.7	54.8
Deep sequential(C)	78.0	77.0	54.7	54.3

Table 7. Results of discourse parsing on multiparty dialogs (F_1 -score). Deep sequential (C) means combine the training set of STAC and Molweni as the training set and test the model respectively.

▼ CASE STUDY

Dialogue 3

nuked: ok likwidoxigen ill reboot and let you know how it goes
likwidoxigen: who makes the printers ? and they wokked before yets ?
nike: yes they worked excellently on dapper . they are two hp deskjets
nbx909: does n't give me the address
likwidoxigen: and they just dont 'print properly ?
likwidoxigen: ok let me keep poking
nbx909: i know but it 's a ups (battery backup) device would it be under sda ?
nuked: i used kde 's add printer wizard , and only samba printers are allowed
likwidoxigen: i 'd assume so . it still has to access the device
likwidoxigen: damn do any usb device work ?

U_1
 U_2
 U_3
 U_4
 U_5
 U_6
 U_7
 U_8
 U_9
 U_{10}

Q1: Who does ask for the address?
Gold answer: *nbx909*
BERT-wwm answer: *likwidoxigen*

Q2: how are printers working?
Gold answer: NA.
BERT-wwm answer: they worked excellently on dapper.

(a)

(b)

Figure 3. Dialogue3. (a) A real example from Molweni dataset with three speakers and ten utterances. (b) Two questions for Dialog 3 and the pridedicted answers of BERT-wwm model.