

NOVAGRAPHS FSA CORPUS

- Dialogues in English between users and a task-oriented conversational agent
- Interactions revolving around the description of Finite State Automata (FSA)
- Part of the NoVAGraphS project, aimed at developing tools for an easier access to educational content, in particular for Visually Impaired People (VIP) in STEM





Data obtained from interactions between users and a rule-based DS providing information on one of two FSA used as examples



Collected in Spring 2023 (Phase 1) + August-September 2023 (Phase 2) Users interacted with the system via a web interface compliant with the Web Content Accessibility Guidelines (WCAG) 2.1 (VIP used a keyboard and a screen reader)

http://www.integr-abile.unito.it/progetto-novagraphs/

3 ANNOTATION DESIGN & RESULTS

In collaboration with:

DIALOGUE ACTS: ISO 24617-2 Standard (Bunt et al., 2017, 2020).

Tagset: 5 dimensions (Task - Ta:, Auto-feedback -Auto:, Turn Management - TuM:, Dialogue Structuring - DS:, Own Communication Management – OCM:) + 6 communicative functions (setQuestion, request, propositionalQuestion, checkQuestion, suggest, answer, autoNegative)

Format: DiAML-TabSW (tabular format, instead of XML); multi-label tagging allowed

LABELS

ERRORS: Adaptation of error annotation scheme from Sanguinetti et al. (2020). Errors intended as violations of gricean maxims (excl. Quality)

Tagset: Quantity: Lack Excess of information; Relation: Ignoring question/feedback, Repetition; Topic change, Straight wrong response, Off topic (new). Manner: Indirect response, Non-understandable, Grammatical error. Generic: Non-cooperativity Format: Tabular; multi-label tagging allowed

Participants' demographics:

- ✤ 6 VIP, 26 non-VIP
- ✤ All non-native English speakers
- ✤ Majority within the age range 25-34, a BSc degree and a background knowledge of FSA, male
- All participated on a voluntary basis



	VIP	Non-VIP	A
# Dialogues	6	26	3
# Turns	194	512	70
Turns/Dialogue	32.33	19.69	22.0
Tokens/User's turn	3.87	5.91	5.5

Basic statistics of the collected data

4 EXPERIMENTS

Users' DA Annotation with DIET classifier

We addressed the task of recognizing users' DAs structuring this process as an intent detection task with DIET, the default classifier provided with RASA

	BoW	BERTemb
Ta:setQuestion	0.9852	0.9963
Ta:request	0.9645	0.9655
Ta:propositionalQuestion	0.9293	0.9583

Pipelines:

- Stag-of-word representation of character ngrams (1 to 4), already set in the default configuration
- BERT pre-trained embeddings
- In both configurations:
- ✤ 100 epochs, 5-fold cross-validation, 3 runs; selection of the three most frequent users' DAs due to insufficient data instances for training

Corpus metadata:

Conversation ID Turn number Participant (User|System) ✤ Text ✤ VIP (yes|no) ✤ No. Tokens FSA ID



Label distribution & inter-annotator agreement

				10 turne to encotate teal, divided into ture main stance
Dialogue Acts		Repetition	42.75	TO lums to annotate; task divided into two main steps:
Users Ta:sotOuestion	51.54	Grammatical error	20.24 13.74	1) dimension; 2) communicative function.
Tatroquest	01.04 07.45	Non-understandable	15.74	Same prompt run three times and independently on each
Ta:propositionalQuestion	16.25		3.82	Came prompt full times and independently on caen
Ta:checkQuestion	1.96	Lack of information	3.05	turn
OCM:selfCorrection	1.40	Non-cooperativity	2.00	For each turn the maximum attainable score is computed
DS:opening	0.56	Ill-formed	1.53	FUI Each luth the maximum attainable score is computed,
TuM:turnAccept	0.56		1.00	considering 1 point per each gold DA (0.5 points per
SOM:initGreeting	0.28	Topio obongo	62.20	aubtaalu dimanajan and aammunjaativa function)
DS		Straight wrong recooner	14.66	sublask: dimension and communicative function).
Ta:answer	48.14	Straight wrong response	14.00	Assigned 1 point per complete, correct DA, 0.5 if only one
AutoF:autoNegative	24.79	Evenes of information	7.00	
DS.enddoet	24.79	Excess of information	7.33	subtask is correct, and U if the annotation is incorrect or if
DO.Suggest	0.07	look of information	4 7 4	
Distribution (in %) of s	2.27 ingle DAs in	Lack of information Ignoring question/feedback Distribution (in %	$\frac{4.71}{0.52}$	the model failed to interpret the prompt correctly.
Ta:suggest Distribution (in %) of s users' and DS turns	ingle DAs in	Lack of information Ignoring question/feedback Distribution (in % errors in users' ar (calculated over to errors)	4.71 0.52 b) of single nd DS turns the total of	the model failed to interpret the prompt correctly.
Ta:suggest Distribution (in %) of s users' and DS turns	2.27 ingle DAs in Cohen's	Lack of information Ignoring question/feedback Distribution (in % errors in users' ar (calculated over t errors) k	4.71 0.52 6) of single nd DS turns the total of	the model failed to interpret the prompt correctly. DATA AVAILABILITY Corpus available for research purposes by filling in a form a resource repository (see OR code)
Distribution (in %) of s users' and DS turns	2.27 ingle DAs in Cohen's Phase	Lack of information Ignoring question/feedback Distribution (in % errors in users' ar (calculated over t errors) k e 1 Phase 2	4.71 0.52 b) of single nd DS turns the total of	the model failed to interpret the prompt correctly. DATA AVAILABILITY Corpus available for research purposes by filling in a form a resource repository (see QR code)
Distribution (in %) of s users' and DS turns	2.27 ingle DAs in Cohen's Phase 0.74	Lack of information Ignoring question/feedback Distribution (in % errors in users' ar (calculated over to errors) k e 1 Phase 2 0.96	4.71 0.52 6) of single nd DS turns the total of	the model failed to interpret the prompt correctly. DATA AVAILABILITY Corpus available for research purposes by filling in a form a resource repository (see QR code)
Distribution (in %) of s users' and DS turns DAs • Communicative Fun	2.27 ingle DAs in Cohen's Phase 0.74 ction 0.91	Lack of information Ignoring question/feedback Distribution (in % errors in users' ar (calculated over to errors) k e 1 Phase 2 0.96 0.96	4.71 0.52 6) of single nd DS turns the total of	the model failed to interpret the prompt correctly. DATA AVAILABILITY Corpus available for research purposes by filling in a form a resource repository (see QR code) Additional data available along with the corpus:
Distribution (in %) of s users' and DS turns DAs • Communicative Fun • Dimension	2.27 ingle DAs in Cohen's Phase 0.74 ction 0.91 0.56	Lack of information Ignoring question/feedback Distribution (in % errors in users' ar (calculated over f errors) k e 1 Phase 2 0.96 0.96 0.96	4.71 0.52 6) of single nd DS turns the total of	the model failed to interpret the prompt correctly. DATA AVAILABILITY Corpus available for research purposes by filling in a form a resource repository (see QR code) Additional data available along with the corpus: •Two PNG files with the graphical representation of the auto

				Users		
	Dialogue Acts			Repetition	42.75	10 turns to annotate; task divide
Users		I	Ignoring question/feedback	28.24	1) dimonsion: 2) communicativo fur	
-	Ta:setQuestion	51.54	(Grammatical error	13.74	
	Ta:request	27.45	1	Non-understandable	4.58	Same prompt run three times and
	Ta:propositionalQuestion	16.25	(Off-topic	3,82	turn
	Ta:checkQuestion	1.96	l	Lack of information	3.05	
	OCM:selfCorrection	1.40	1	Non-cooperativity	2.29	For each turn the maximum attair
	DS:opening TuM:turnAccont	0.56		III-formed	1.53	concidering 1 point per each a
	SOM initGreeting	0.56		DS		considering i point per each g
-	DS	0.20		Topic change	62.30	subtask: dimension and communication
-	Ta:answer	48.14	5	Straight wrong response	14.66	Accianad 1 naint nor complete or
	AutoF:autoNegative	24.79	I	Indirect response	10.47	Assigned i point per complete, co
	DOussianast	24 70		Excess of information	7.33	
	DS:suggest	24.75				subtask is correct, and 0 if the ar
- Distribu	Ta:suggest ution (in %) of si	2.27 2.27	I I As in E	Lack of information Ignoring question/feedback Distribution (in %	4.71 0.52) of single	subtask is correct, and 0 if the ar the model failed to interpret the pro e
Distribu users' a	Ta:suggest ution (in %) of si and DS turns	ingle D	As in E	Lack of information Ignoring question/feedback Distribution (in % errors in users' an calculated over t errors)	4.71 0.52) of single nd DS turns the total of	subtask is correct, and 0 if the another model failed to interpret the prosent of DATA AVAILABILITY
- Distribu users' a	Ta:suggest ution (in %) of si and DS turns	ingle D	As in E As in E ((ohen's k	Lack of information Ignoring question/feedback Distribution (in % errors in users' an calculated over t errors)	4.71 0.52) of single nd DS turns he total o	subtask is correct, and 0 if the another model failed to interpret the pro of DATA AVAILABILITY Corpus available for research purp
Distribu users' a	Ta:suggest ution (in %) of si and DS turns	ingle D	OAs in [0As in [6 (6 0hen's k Phase 1	Lack of information Ignoring question/feedback Distribution (in % errors in users' an calculated over t errors) Phase 2	4.71 0.52) of single nd DS turns the total o	subtask is correct, and 0 if the ar the model failed to interpret the pro DATA AVAILABILITY Corpus available for research purp resource repository (see QR code)
Distribu users' a	Ta:suggest ution (in %) of si and DS turns	ingle D	OAs in [6 6 6 6 6 6 6 7 6 7 6 7 7 7 7 7 7 7 7	Lack of information Ignoring question/feedback Distribution (in % errors in users' an calculated over t errors) Phase 2 0.96	4.71 0.52) of single nd DS turns the total o	subtask is correct, and 0 if the art the model failed to interpret the pro DATA AVAILABILITY Corpus available for research purp resource repository (see QR code
Distribu users' a DAs • Cor	Ta:suggest ution (in %) of si and DS turns mmunicative Fund	ingle D	As in [As in [6 (4 6 0hen's k Phase 1 0.74 0.91	Lack of information Ignoring question/feedback Distribution (in % errors in users' an calculated over t errors) Phase 2 0.96 0.96	4.71 0.52) of single nd DS turns the total o	subtask is correct, and 0 if the ar the model failed to interpret the pro DATA AVAILABILITY Corpus available for research purp resource repository (see QR code) Additional data available along w
Distribu users' a DAs • Cor • Dim	mmunicative Fund	ingle D	As in E OAs in E (4 0 0 0 0 0.74 0.91 0.56	Lack of information Ignoring question/feedback Distribution (in % errors in users' an calculated over t errors) Phase 2 0.96 0.96 0.96 0.96	4.71 0.52) of single nd DS turns the total o	subtask is correct, and 0 if the ar the model failed to interpret the pro DATA AVAILABILITY Corpus available for research purp resource repository (see QR code) Additional data available along w •Two PNG files with the graphical

macro-F1	0.9597	0.9734

purposes using DIET

Annotating DAs with LLMs

LLMs *zero-shot* performance in Experiments on annotating DAs

Prompt building:

All prompts comprised the following points:

task definition (i.e. annotating DAs in turns)

Models tested: ChatGPT 3.5 (via web interface) ✤ Llama 2

✤ Tk-Instruct

context (i.e. the dialogue between a student and a DS programmed to answer about a specific FSA)

A annotation constraints (e.g. one or more DAs can be assigned to each turn, additional co-text)

presentation of the label set

DAs decomposition into dimension and communicative function

Propositional questions to guide the selection of the most appropriate dimension and communicative function

the turn to annotate

Prompt	ChatGPT	Llama 2	Tk
1	0.33/1	0.00/1	0.00/1
2	0.33/1	0.00/1	0.50/1
3	1.00/1	0.00/1	0.50/1
4	0.00/1	0.17/1	0.00/1
5	0.50/1	0.00/1	0.00/1
6	0.17/1	0.00/1	0.50/1
7	0.83/2	0.00/2	0.50/2
8	1.00/1	0.00/1	0.50/1
9	1.00/2	0.00/2	0.50/2
10	1.17/3	1.00/3	0.50/3

0 00/4 40 0 44/4 40 0 0 0 0 14 40

average	0.63/1.40	0.11/1.40	0.35/1.40
SD	0.41/0.70	0.31/0.70	0.24/0.70
success	45.24%	8.33%	25.00%

poses by filling in a form available in the

vith the corpus: representation of the automata ate tables of the automata



Inter-Annotator Agreement

References

• Harry Bunt, Volha Petukhova, David Traum, and Jan Alexandersson. 2017. Dialogue act annotation with the iso 24617-2 standard. In Deborah A. Dahl, editor, Multimodal Interaction with W3C Standards: Toward Natural User Interfaces to Everything, pages 109–135. Springer. • Harry Bunt, Volha Petukhova, Emer Gilmartin, Catherine Pelachaud, Alex Fang, Simon Keizer, and Laurent Prévot. 2020. The iso standard for dialogue act annotation. In 12th Edition of its Language Resources and Evaluation Conference (LREC 2020), pages 549–558. ELRA

• Manuela Sanguinetti, Alessandro Mazzei, Viviana Patti, Marco Scalerandi, Dario Mana, and Rossana Simeoni. 2020. Annotating errors and emotions in human-chatbot interactions in Italian. In *Proceedings of the 14th Linguistic Annotation Workshop*, pages 148–159, Barcelona, Spain. Association for Computational Linguistics