

Reference Resolution and New Entities in Exploratory Data Visualization: From Controlled to Unconstrained Interactions with a Conversational Assistant

Abari Bhattacharya^{*1}, Abhinav Kumar^{*1}, Barbara Di Eugenio¹,
Roderick Tabalba², Jillian Aurisano³, Veronica Grosso¹,
Andrew Johnson¹, Jason Leigh², and Moira Zellner⁴

¹University of Illinois Chicago

{abhattach62, akumar34, bdieugen, vgross3, ajohnson}@uic.edu

²University of Hawaii at Manoa {tabalbar, leighj}@hawaii.edu

³University of Cincinnati jillian.aurisano@uc.edu

⁴Northeastern University m.zellner@northeastern.edu

Abstract

In the context of data visualization, as in other grounded settings, referents are created by the task the agents engage in and are salient because they belong to the shared physical setting. Our focus is on resolving references to visualizations on large displays; crucially, reference resolution is directly involved in the process of creating new entities, namely new visualizations. First, we developed a reference resolution model for a conversational assistant. We trained the assistant on controlled dialogues for data visualizations involving a single user. Second, we ported the conversational assistant including its reference resolution model to a different domain, supporting two users collaborating on a data exploration task. We explore how the new setting affects reference detection and resolution; we compare the performance in the controlled vs unconstrained setting, and discuss the general lessons that we draw from this adaptation.

1 Introduction

Conversation is understood in context. When the world, whether real or simulated, can change because of the user's actions, new entities are created by the processes that change the world itself: then, reference resolution, which links what the user refers to with objects in the world, is crucial for a dialogue system to effectively respond to the user, including by creating new entities.

Our overall research program aims to develop and deploy flexible conversational assistants to support users, whether causal or professional, and whether alone or in teams, explore data via visualizations on large screen displays - large screen

displays better support exploration and collaboration (Andrews et al., 2011; Rupprecht et al., 2019; Lischke et al., 2020). In this paper, we focus on new entity establishment via reference in such contexts. We start from the corpus *Chicago-Crime-Vis* we collected a few years back (Kumar et al., 2016, 2017) in which a user exploring crime data in Chicago interacts with a Visualization Expert (VE) whom they know to be a person generating visualizations on the screen remotely from a separate room. On the basis of *Chicago-Crime-Vis*, we designed and developed a version of our assistant which was called *Articulate2* (Aurisano et al., 2016; Kumar et al., 2020)¹. We will report the performance of *Articulate2* on reference resolution, and especially reference establishment, with respect to the transcribed and annotated *Chicago-Crime-Vis* corpus, evaluated in an offline manner. The second part of our paper discusses the challenges that arose when we ported *Articulate2* to a new setting: two collaborators work together to assess COVID policies given geographic and demographic features of the data, and interact exclusively with the deployed *Articulate+* (see Figure 1). We will illustrate the many issues which degrade performance, from speech processing errors, to the adaptation of models to new domains, to the inherently more complex setting in which the assistant is now behaving like an overhearer of somebody else's conversations. For clarity, we will refer to *Articulate2* in the city crime domain as *Art-City-Asst*, and to *Articulate+* in the COVID domain, as *Art-COVID-Asst*.

A disclaimer before we proceed: the purpose

¹The first interface we developed in this space was called *Articulate* (Sun et al., 2010).

*Co-first authors

of this work was to adapt a previously developed conversational assistant and to evaluate it in a more unconstrained setting. We do not believe in chasing after the latest shiny approach, including ChatGPT², and undertake a potentially infinite loop of changes which would never bring us to real user studies. Additionally, we strongly believe in ecologically valid data, such as our *Chicago-Crime-Vis* data. This data is by nature small, in fact tiny as compared to most current datasets. We will return to these issues in the Conclusions.



Figure 1: User setting for COVID data exploration, with two collaborators

2 Related Work

2.1 Conversational assistants for data visualization

Earlier work on conversational assistants for data visualization include (Cox et al., 2001), which established the benefits of using NL to generate visualizations for exploratory data analysis. In the ensuing 20 years, several such systems emerged in this area, see (Shen et al., 2023) for a systematic survey: e.g., DataTone (Gao et al., 2015a), FlowSense (Yu and Silva, 2020), Eviza (Setlur et al., 2016a) and DT2VIZ (Jiang et al., 2021). Lately, Large Language Models (LLMs) have started being integrated into visualization tools (e.g., PandasAI from OpenAI³), but not as part of a conversational assistant that keeps track of dialogue history.

Our previous work - the Articulate assistant series. Our research program started more than 10 years ago with *Articulate*, one of the first conversational assistants for creating data visualizations (Sun et al., 2010), as also noted by (Shen et al.,

2023). The first *Articulate* would only respond to individual commands, but even so, users were 12 times faster when using *Articulate* to generate a chart in comparison to a spreadsheet program (Microsoft Excel). Still, the commands that *Articulate* would answer to were not grounded in actual human data; hence, we collected the *Chicago-Crime-Vis* corpus (Aurisano et al., 2015; Kumar et al., 2016, 2017) that informed a new prototype, *Articulate2*, a multimodal system that could support speech commands and gestures to facilitate data exploration tasks (Kumar et al., 2020; Kumar, 2022); and whose reference resolution component we are discussing in this paper. Subsequently, we ported *Articulate2* to the COVID domain, dubbed it *Articulate+* and developed two versions of the NLI: *Articulate+-PE* and *Articulate+-DM*. *Articulate+-PE* (Tabalba et al., 2023, 2022), was developed independently (from scratch), and works by identifying database properties or attributes mentioned directly or indirectly in the utterances. To identify the chart types given the utterance, it uses a Chart Classifier Neural Network trained on a small dataset of utterances from a preliminary user study using NLP.js library⁴. However it lacks dialogue management as well as reference resolution. The other version, *Articulate+-DM*, is *Articulate2* ported to the COVID domain. To reiterate then, in this paper we discuss the evaluation of *Articulate2*, in its incarnation as *Art-City-Asst* evaluated offline on the *Chicago-Crime-Vis* data (Section 5), and in its second incarnation in the COVID domain as *Art-COVID-Asst* evaluated in an actual user study (Section 6).

2.2 Co-Reference Resolution

This field is as old and as vast as NLP; here we focus on its applications to visualization, which are hindered by several limitations: e.g., only referents to objects within the current visualization are handled (Sun et al., 2010; Gao et al., 2015b; Narechania et al., 2020), or only referents for follow-up queries on a current visualization are tracked (Reithinger et al., 2005; Setlur et al., 2016b; Hoque et al., 2017; Srinivasan and Stasko, 2017). As (Shen et al., 2023) concludes, "existing [approaches] mostly leverage NLP toolkits to perform co-reference resolution. Although useful, they lack detailed modeling of visualization elements" or, we would add, of what has transpired earlier in the

²<https://openai.com/product/chatgpt>

³<https://www.kdnuggets.com/2023/05/pandas-ai-generative-ai-python-library.html>

⁴<https://github.com/axa-group/nlp.js/>

dialogue. In contrast to this, we focus on reference resolution within an environment in which visualizations are dynamically added to and removed from the screen, and can subsequently be referred to. This requires *accommodating context change*, a notion first introduced by (Webber and Baldwin, 1992) in their discussion of new entities that are the results of physical processes as in cooking (e.g., *the dough* resulting from *mixing flour, butter and water*). In the 30 years since, not much work has been done on how to accommodate the creation of new entities⁵ (see (Wilson et al., 2016) for documents and (Li and Boyer, 2016) for tutoring dialogues about programming), and none in the visualization domain. Note we do not focus on multimodal reference resolution, another vast area (Navaretta, 2011; Qu and Chai, 2008; Eisenstein and Davis, 2006; Prasov and Chai, 2008; Iida et al., 2011; Kim et al., 2017; Sluÿters et al., 2022), even if we will briefly touch on deictic gestures in Section 3.

3 Controlled Dataset: Chicago-Crime-Vis

Our *Chicago-Crime-Vis* corpus comprises multimodal interaction for 16 subjects that explored public crime data in our city to better deploy police officers.⁶ As noted, they spoke with a human VE who remotely created visualizations on a large screen, was not visible and did not speak back. The corpus contains 3.2K utterances. Since the user was encouraged to reason out loud about the patterns discovered via visualization, conversational turns often start with *think aloud*, followed by what we call an *actionable request* (AR) for the VE.

Using ANVIL (Kipp, 2001, 2014), we annotated 449 CARs (*contextual actionable requests*), covering 1545 utterances: a CAR consists of *setup*, i.e. think aloud prior to the AR (up to and including utterances that mention data attributes, if any); the AR; and the *conclusion*, the think aloud subsequent to the AR (also based on data-attribute mentions). While each AR is just one utterance, each of set-up and conclusion may include more than one —on average, 1.8 and 2 respectively. (See Table 1, *Chicago-Crime-Vis*(H) column for the distribution of set-ups and ARs annotated in the dataset). Fig-

⁵Work in formal pragmatics that models extra-linguistic context exists - e.g. see (Stojnic et al., 2013; Hunter, 2014), but as far as we know, it has not been used to model references in actual physical contexts.

⁶We acknowledge that this task may be fraught in the era of Black Lives Matter in the United States. This data was collected prior to 2020, when the current awakening as concerns policing and racism surfaced to public consciousness.

	Chicago-Crime-Vis (H)	COVID (A)	COVID (T)
Set-up	218	73	149
AR	449	1296	2563

Table 1: Total count of Set-ups and ARs in the 3 user studies —H: Human; A: Automatic; T: Transcript

ure 2 shows two CARs from our corpus, which we will use as our running example.

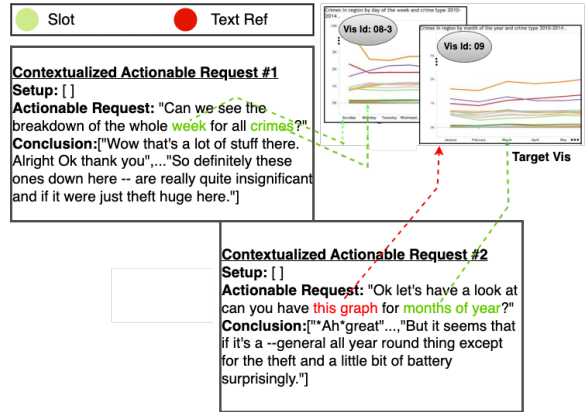


Figure 2: Excerpt comprising two CARs; references shown in red and slot fillers in green. In CAR #1 visualization "08 – 3" is specified via temporal axis DAY associated with slot filler "week" and similarly CRIME for "crimes". CAR #2 creates "09" substituting temporal axis DAY in "08 – 3" with MONTH, associated with slot filler "month of year". The identifiers are internal to the system but not visible to the users.

Each AR is annotated for user intent with one of 8 Dialogue Acts (DA) labels (with excellent intercoder agreement on the 8-way annotation, $k = 0.74$), including: WINMGMT for window management operations, e.g., closing, or minimizing; CREATEVIS for creating a new visualization from scratch; MODIFYVIS for creating a new visualization based on an existing one. The transcribed corpus is publicly available⁷, and so is an augmented dataset built to alleviate data scarcity, comprising a 10-fold increase to 160 subjects covering approximately 15K utterances obtained via delexicalization and paraphrasing.

Referring Expression Annotation. We annotated both *text* (NPs) and *gestural* references to visualizations. Hand gestures were coded with various labels (e.g., the kind of gesture, the objects pointed to on the screen, and so on); approximately

⁷<https://github.com/uic-nlp-lab/Chicago-crime-vis-corpus>

Category	Setup	AR
Overall	19	109
Single Referents	18	86
Single Targets	14	66

Table 2: *Chicago-Crime-Vis* text reference distribution

a third were identified as referential when they co-occur with text references. We labeled a total of 294 references in the 449 CAR’s, of which 176 textual, and 118 gesture. We obtained an excellent intercoder agreement of $\kappa = 0.85$ with 2 judges on the full interaction from one subject. Given lack of space, and because in our unconstrained setting gestures were not addressed, we will not discuss gestures further. Table 2 shows the text reference distribution where within the 176 text references (of which 19 appear in set-up, 109 in AR, and 58 in conclusions). We also annotated 680 phrases as slot fillers corresponding to data attributes (i.e., *slots*) in our knowledge ontology (KO). The KO was semi-automatically constructed via external sources such as our city portal, augmented with synsets extracted from Wordnet⁸ and Babelnet⁹; it comprises 3.5K total terms categorized into 11 parent types such as CRIME TYPE, NEIGHBORHOOD, TIME etc, of which about half are common nouns and about half proper nouns pertaining to Chicago.

4 Co-Reference: Detection, Resolution, and New Entity Establishment

We briefly discuss the NLP engine (in the context of the full conversational assistant, see Figure 3), focusing on its reference resolution component - full details on the NLP engine can be found in (Kumar et al., 2020; Kumar, 2022). The NLU pipeline relies on an information state architecture with dialogue state tracking. After speech recognition (please see below for further discussion), traditional parsing and semantic role labeling are performed, and then a semantic frame is computed (see below). The dialogue management module is responsible for: classifying the intent of the user as one of the 8 DAs mentioned in Section 3; performing reference resolution; and updating and maintaining the dialogue history (DH). The NLP engine transforms the user request (when appropriate) into an SQL query; and in a visualization specification

⁸<https://wordnet.princeton.edu/>

⁹<https://babelnet.org/>

that is passed to Vega-Lite¹⁰, a separate visualization interface software, to create a visualization of the data returned by the SQL query and add it to the display.

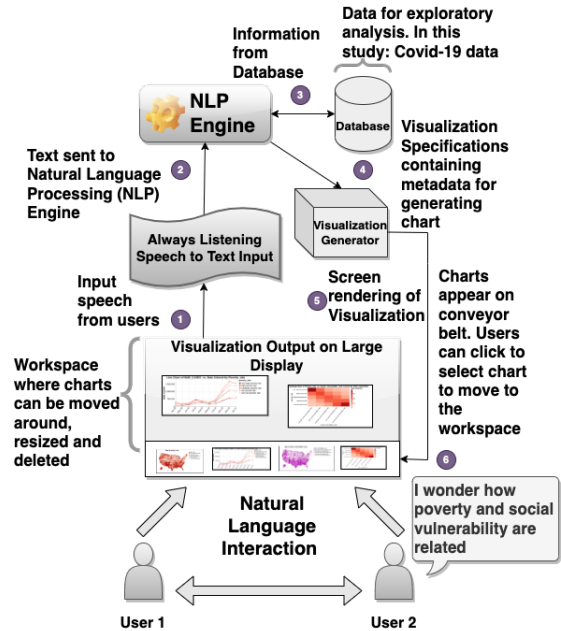


Figure 3: The Conversational Assistant—in its COVID incarnation, with two collaborators. The annotated arrows denote the workflow of the architecture, the numbers signify the order of events when the users interact with the conversational assistant.

4.1 Semantic Frame Construction

Each time a visualization is mentioned in the dialogue (whether it refers to a previous one or not) our model looks for slots in the request to form its semantic frame. We find phrases that are in close proximity in the embedding vector space to terms in the KO, by using a domain targeted word embedding model (WEM)¹¹. Subsequently the candidate words are pruned based on linguistic patterns using the SpaCy¹² dependency parse of the entire utterance to form the final list of slot fillers. For example in the AR in CAR #2 in Figure 2, the prepositional phrase “for months of year” contains “month” and “year”, both of which are known as temporal slots in KO. Here, the terms are merged to form “months of year”, and mapped to the parent slot MONTH - see User Action (1) in Figure 4.

¹⁰<https://vega.github.io/vega-lite/>

¹¹100-dimensional continuous bag-of-words model trained on 5GB of online articles and wikipedia pages related to crime.

¹²<http://spacy.io>

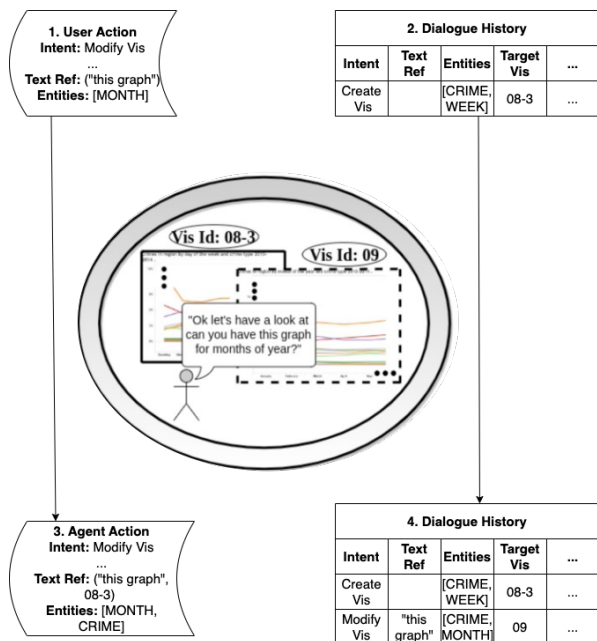


Figure 4: The user (inside the circle) currently has visualization "08-3" on the screen and is asking to construct a new visualization "09"(in dashed lines since it is being built - visualization identifiers are internal to the system but not visible to the users). Reference resolution operates in four stages. NLU creates user action (1); the DM uses DH (2) to create agent action (3); finally the state tracker updates DH (4).

4.2 Dialogue Manager (DM)

The DM executes a dialogue policy which aside from making back-end decisions such as forming an SQL query for data retrieval, also seeks to populate unknown frame attribute values - semantic frames are constructed in response to either a CREATEVIS or a MODIFYVIS DA, and in the case of MODIFYVIS, reference resolution may be used to fill some of those unknown values. When the semantic frame is complete, the state tracker adds it as a new entry to the DH while the system also outputs a json object (which we call a *visualization specification*) that instructs Vega-Lite2 to accordingly update the screen.

For example, in Figure 4, in CAR 1, AR #1 "Can we see the breakdown of the whole week for all crimes"? has resulted in updating both the DH and the screen with a new linechart (the new visualization "08-3"). After AR #1, the DH contains a single entry for "08-3" and its specifications in a frame-slot format, including: the user intent (CREATEVIS), the type of plot, and its semantic frame in terms of attributes that were mentioned (*crime, week*) - see Dialogue History 2 in Figure 4. Note

that IDs like "08-3" are for internal reference, and not shown to the user, but are included in Figures 2 and 4 for ease of exposition.

When AR #2 is processed and a MODIFYVIS DA is recognized, a new frame is created (see Agent Action 3 in Figure 4); while user intent (MODIFYVIS) and some slots (MONTH) are filled, others are left empty either because of under-specification by the user (e.g., axes labels, plot type, and so on) or they require additional processing by the DM; in this particular case, the previous visualization "08-3" will be found as the referent for *this graph* and both CRIME will be added as an additional slot, and the plot type will be inferred to be line chart (see below) - see Dialogue History 4 in Figure 4.

Next, we describe reference detection and resolution.

4.2.1 Reference Detection

We trained a sequence tagging model to detect text references (DTR). The model predicts tags using the standard IOB2 format (i.e., "B-REF"/"I-REF"/"O-REF" for beginning of / inside / outside text reference respectively). We trained a simple CRF model that uses POS tags as features, and two baseline models, BiLSTM-CRF and BERT-CRF. Further, to remedy data insufficiency - there are only 176 text references appearing across 449 CARs in the corpus, we investigated Sequential Transfer Task Learners and Multi Task Learners, in both cases, as applied to BiLSTM-CRF and BERT-CRF. As transfer or additional task, we use a NER task based on our augmented dataset, which is also automatically labelled for 23 NER tags, based on the B/I/O scheme: the "B" and "I" tag for each of the 11 parent slots in the KO (e.g., B-visualization, I-visualization) plus "O" tag (the slot names are known because they are manually labelled in the 449 CARs and delexicalization maintains their type).

4.2.2 Reference Resolution

To understand to which visualization the current referring expression refers, we use heuristics based on recency and similarity. The slot fillers from the frame of the current referring expression and from the candidate visualizations in the DH are transformed into *visualization vectors*, ie, they are projected onto an embedding space along 11 dimensions, corresponding to the 11 slots in the KO, using the WEM mentioned earlier. Before compar-

ing the two visualization vectors, a recency factor is applied. If n represents the total entries in the DH, then the visualization vectors of the most recent $\frac{n}{2}$ entries in the DH are associated with a multiplicative factor of 1.0 signifying that they are equally preferred. The latter $\frac{n}{2}$ entries in the DH however are associated with a linear decrease by a factor of $\frac{1}{n}$. Finally, cosine similarity is used to score each visualization in the DH relative to the referring expression and the visualization with the highest score is selected, as long as it exceeds a cut-off of 0.40 (established empirically).

For example in Figure 4, the DH contains only an entry for "08 – 3" (other earlier visualizations must have been closed and are not relevant any more). Since the cosine similarity score between "08 – 3" and the current semantic structure exceeds 0.4, "08 – 3" is chosen as the referent for *this graph*.

4.2.3 New Entity Establishment

Once the referent of the specific referring expression has been established, a new visualization ("09") is constructed using the referent's frame representation to infer missing information ("08–03"). Explicit information in the current request is used to replace identical slots: e.g. MONTH, which was used to resolve the referring expression via WE embedding and cosine similarity among semantic structures, replaces WEEK as the temporal axis in "09". Information that is unspecified in the request but present in the referred-to visualization is imported to establish the new visualization; in this particular case, CRIME is added to the slot list because "08 – 3" of the previous request includes it. Finally, to generate the new visualization corresponding to a referring expression, the chart type (heat map, line chart, or bar graph) also needs to be inferred; it is simply copied from the referent, resulting in the new linechart "09" being added to the screen, and the updated entry being added to DH (#4 in Figure 4).

5 Constrained Evaluation on Chicago-Crime-Vis

The results we present now were obtained by manually evaluating the pipeline, which was run on the transcribed *Chicago-Crime-Vis* data in an offline manner: hence, we did not have to contend with speech errors, or with error propagation, since for every utterance, the DH up to that point was reset to a correct state if necessary. Currently, our model

focuses on references occurring in *setup* and *AR* for detection, and in *AR* only for evaluation of semantic frame correctness. Additionally, we focus on single referents and single targets: e.g. in "*Can you bring up the graph behind the River North one?*" the user refers to two visualizations; whereas "*well I would like to see battery by day of week, battery by month, and battery by year.*" results in 3 new corresponding visualizations. However, our model only adds one of these visualizations to the dialogue history (DH) as part of the evaluation. Table 2 presents text reference counts only for *setup* and *ARs* (hence, excluding 58 references in *conclusion*). Single referents account for about 94.7% of references in *setup* and for about 80% of those in *ARs*. Finally, when filtering on single targets, we are left with the 80 text references (last row in table) on which we will focus.

5.1 Detection

Notwithstanding the lack of training data, the CRF performed the best, achieving a 61.2% F1 on the B-REF, I-REF, O-REF task. This is statistically significantly better than any other models (the next best is Multitask BERT-CRF with F1= 43.5%). Hence, the CRF model is used in the subsequent steps in the pipeline. The five-fold cross validation accuracy of this CRF model on the *Chicago-Crime-Vis* data is shown in Table 3.

5.2 Resolution

Accuracy on resolving text references for varying *WINDOW* sizes is shown in Table 4. If one only takes into account the visualization introduced by the preceding *AR* (recall that we currently don't deal with multiple references), accuracy is 85.3% for *set-up* and 74.4% for *AR*. Interestingly, in the *Chicago-Crime-Vis* corpus, users also refer to the most recent visualization over 75% of the time. However, when we provide unlimited window size (∞ means all referent visualization candidates are eligible), resolution of references in *ARs* decreases; this suggests our linear decay function may need further tuning to better model the user preference behavior.

5.3 Semantic Frame Accuracy

We report the performance of semantic structure construction as concerns *CREATEVIS* and *MODIFYVIS* *AR*'s. Our model achieved a slot accuracy metric (Takanobu et al., 2020) of 66.2% for semantic slots: this concerns the specification of the slots

	Chicago-Crime-Vis	COVID (A)	COVID (T)
Set-up	60.0	50.0	33.3
AR	55.0	25.0	45.8

Table 3: Evaluation of reference detection model. Chicago-Crime-Vis: five-fold cross validation accuracy calculated on Single Targets of Table 2; COVID (A): Accuracy in real-time user study; COVID (T): Accuracy on correct transcripts of real-time user study. COVID (A) and COVID (T) evaluated on a significant sample size

	Setup Window		AR Window	
	1	∞	1	∞
Chicago-Crime-Vis	85.3	85.3	74.4	68.3
COVID (T)	-	-	36.3	54.0

Table 4: Resolution accuracy for varying window sizes. COVID (T) evaluated on a significant sample size

of the *Visualization Frame (VH)* in the DH, and includes slots that were explicit in the utterance, and those that were inferred. Given the example in Figure 2, for "08 – 3" the two slot values are "crime" and "week", and for "09" "month" (explicit) and "crime", inferred via reference resolution. Table 5 reports the number of VFs for which a certain percentage of slots has been correctly recognized, by quartile. The 100% quartile is equivalent to the *Joint Goal Accuracy (JGA)* metric used in some of the Dialogue State Tracking challenges, which *compares the predicted dialog states to the ground truth at each dialog turn, and the output is considered correct if and only if all the predicted values exactly match the ground truth* (Takanobu et al., 2020). For the *Chicago-Crime-Vis*, these were manually annotated when annotating for references, and the results are computed by evaluating the resolution pipeline turn by turn, with the gold-standard DH up to the previous turn: in 131 of those (55%), all slots were correctly recognized; in 83% of these VFs, at least 75% of the slots were correct; only in 17 (7%) of these 238 VFs, no slots were correctly recognized. Beyond Joint Goal Accuracy, we report partial accuracy to provide a more nuanced analysis of the assistant’s performance, which cannot be simply measured in a binary "Correct/Incorrect" fashion: in an dialogue based application for data exploration like ours, a partially recognized visualization frame can generate charts which may help the users move forward. Papers exploring similar views are Selfridge et al. (2011) and Schlangen

	0%	25%	50%	75%	100% (JGA)	Total VF
Chicago-Crime-Vis	17	5	19	66	131	238
COVID (A)	22	1	25	8	66	122
COVID (T)	23	4	25	15	75	142

Table 5: Distribution of *Visualization Frames* wrt % correct slots. COVID (A) and COVID (T) evaluated on a significant sample size

et al. (2009), where partial speech recognition and reference resolution were found to be beneficial for dialogue systems that react satisfactorily to the user.

6 Unconstrained setting: User studies in a COVID domain

A realistic evaluation of the NLP Architecture was conducted through user studies: pairs of participants interact with the conversational assistant (*Art-COVID-Assst*) to perform two open-ended exploratory data analysis tasks, concerning which factors may affect COVID mitigation strategies, such as access to doctors or elderly population. Overall, 15 groups of 2 participants, performed the two tasks in a specified sequence, within a time limit of 25 minutes per task. The participants, aged 18+ , were recruited from UIC and were mostly graduate students. With their consent, we audio and video recorded them, and collected logs generated by the back-end code of *Art-COVID-Assst* for analysis purpose. As shown in Figure 1, they are sitting and wearing a mike; also, each has a mouse with which they are able to reposition and click the visualizations on the screen. We encouraged the users to freely interact with each other and with *Art-COVID-Assst*, and we did not provide specific instructions about the tasks, the interface, or the collaboration. The system is designed to "always listen" to the participants, whether or not they are addressing the assistant directly. This is implemented using the Web Speech API¹³.

It was relatively simple to port *Art-City-Assst* to *Art-COVID-Assst* (Figure 3 shows the architecture) and mostly required to update the KO. For the COVID data, we identify 13 semantic slots like "COVID vulnerability rank", "Access to doctors", "Diabetes risk", "Uninsured rate" etc. and the possible values for these slots. As earlier, we enlarged the KO with synonyms for each slot and their values by using Wordnet and Babelnet to generate

¹³<https://wicg.github.io/speech-api/>

these synonyms. The generated KO has a vocabulary of 710 terms. This, as we describe in Section 4 forms the backbone of semantic slot filling and new entity establishment. We keep the same Dialogue Manager as before and use the best Reference Detection model built using the *Chicago-Crime-Vis* corpus, namely, the CRF model. The Reference resolution algorithm also remains the same. Finally for screen rendering of the generated charts, the relevant data obtained from the database is converted to Vega-Lite grammar.

6.1 Findings of the User Study

To evaluate the reference detection and resolution pipeline in this setting, in principle we only need the log of the interactions to assess real-time performance wrt the utterances from the conversations of the participants. However, after we realized that speech recognition errors were a major bottleneck in the real-time study, we conducted additional experiments on the transcripts. These are generated using the Whisper speech recognition model¹⁴ followed by light manual inspection. The corrected transcripts are then fed to the back-end code of the conversational assistant and new logs are generated. We name this version of the user study data as COVID (T) (for *Transcript*), while the real-time logs are named COVID (A) (for *Automatic*).

Since, as we noted earlier, reference detection applies to set-up and ARs, Table 1 shows the distribution of setups and requests in these two versions along with those from the *Chicago-Crime-Vis* corpus. An important difference is that set-ups and ARs for *Chicago-Crime-Vis* were manually annotated, whereas these are the results of automatic recognition for the COVID study (whether A or T). The table shows that there are many more set-ups in the *Chicago-Crime-Vis* data; this difference is significant, as confirmed by $\chi^2 = 489.9511, p < 0.00001$ (with Bonferroni correction). There may be various reasons for this, one being that the classifiers that recognize setup and ARs were trained on the augmented *Chicago-Crime-Vis* corpus and perform worse here to start with. However, it is also possible that in fact, think aloud that feels natural when somebody is by themselves is not in a collaborative situation: a set-up by definition doesn't talk about a data attribute, but we surmise that the two collaborators are more focused on data attributes

¹⁴<https://github.com/openai/whisper> - it became available in September 2022, after our conversational assistant was developed and hence could not be used for the user study.

than on thinking aloud, precisely because they are interacting with another person.

For the purpose of the evaluation, we need to manually verify the results returned by the reference pipeline. Given the size of the data, we obtain two samples, one from COVID (A) (# utterances: 3096) and one from COVID (T) (# utterances: 8440). A significant sample size is computed for both with 95% confidence interval and 5% margin of error. This results in a random sample of 340 (11%) utterances for COVID (A), and of 370 (4.38%) utterances for COVID (T). Subsequently, we use COVID (A) and COVID (T) to refer to these samples of the respective groups, not to the whole group; all evaluation and analysis are done on these samples only.

6.1.1 Reference Detection

Table 3 shows the accuracy of the detected references in Set-up and Request utterances of COVID (A) and COVID (T). As expected, the performance degrades in a real-time user study scenario. Unlike the controlled study setting with one participant, when two people collaborate for an exploratory task, three things happen. First they talk to each other; next, they make requests to the system and finally they draw conclusions. These make reference detection in utterances extremely complex. In the case of COVID (A), we also attribute the lack of accuracy to speech-recognition errors.

6.1.2 Reference Resolution

We limit the evaluation of the reference resolution pipeline to COVID(T) as there were no references resolved during the actual study—DAs of around 44% of those utterances with detected references were misclassified (note that useful visualizations may have been created all the same in response to those specific utterances, but not because a reference resolution was resolved). After conducting a thorough manual inspection of the issue we find the speech recognition errors to be the major roadblock yet again. However using the corrected transcript (COVID (T)) we get a comparatively better performance as shown in Table 4. Since in this study setting, only ARs where references are detected are resolved, we limit our evaluation to ARs only. Contrary to the constrained *Chicago-Crime-Vis* setting, where considering only the previous AR was the better strategy, here limiting window size to 1 results in lower accuracy. We observe that in a more real scenario, especially when two peo-

ple are involved in the conversation, there are more relevant entries in the dialogue history. This may also be due to the nature of the interaction with the large screen: in *Chicago-Crime-Vis*, the user was standing in front of the large display, and often fairly close so that they would in fact mostly focus on only a portion of the display; in the COVID study, the two collaborators were sitting at about 6 ft from the screen (see Figure 1), and hence all visualizations on the screen are more readily available to them.

6.1.3 Semantic Frame Accuracy

For the user study settings of *Art-COVID-Asst*, VFs were recognized for utterances having DAs CRE-ATEVIS and MODIFYVIS. Similar to what we observed in the controlled setting of *Chicago-Crime-Vis* (as described in Section 5.3) in Table 5, more than 50% VFs had all their semantic slots recognized as fully correct in the unconstrained settings with *Art-COVID-Asst*. In fact, we see comparable performances of COVID (A) and COVID (T) across all quartiles. This shows that irrespective of the problematic performance of the speech-to-text algorithm, more than 60% VFs had 75% or more slots correctly filled and more than 80% VFs had at least 50% slots correctly identified. This also explains the reasonable success of the user study that we observed despite the subpar performance of the speech-to-text algorithm. This is attested by questionnaires the users filled. On a 5 point Likert scale, mean scores of 4 and 3 were respectively obtained for usefulness of the charts generated, and for ease of command system use.

7 Conclusions and Future Work

We have presented a reference resolution model for conversational assistants that help user in exploratory data visualization. In particular, the model resolves visualization references in the context of the current interaction, crucially tracking visualizations constantly being added to the screen. The model is central to the creation of new visualizations: visualization features encoded in the DH as slot values, help the model know how to refer to a visualization later on. We have also shown how the initial assistant, *Art-City-Asst*, was ported to a completely different domain. We presented the evaluation of the reference pipeline in both settings, the constrained *Chicago-Crime-Vis* and the "wild" COVID setting, in which two collaborators were exploring COVID data. We are fully aware that

our results are not compared to an external baseline, but we contend that evaluations in grounded settings are important, and do not require creating some artificial baseline or evaluating the pipeline on existing reference resolution datasets.

Not surprisingly, the user evaluation brought several issues to the fore. First, we discovered that the speech API that we had chosen did not work very well (it would have been impossible to change it during the user study even if we had noticed it). Whereas this is unfortunate, we were able to obtain correct transcripts and run a second evaluation. Second, the nature of the interaction and the setting affected the conversations and the results: for example, we found many fewer set-ups in the COVID data, but on the other hand, more references to referents further back in the conversation.

Potential extensions for future work include ways to better model user behavior for referring to more distant visualizations and using sophisticated machine learning approaches in our resolution algorithm to take advantage of the rich visualization feature space in our case. Additionally, in the COVID user study users don't use hand gestures to interact with the screen, however they do use their mouses to click and reposition visualizations, hence bringing multimodality to the fore; not to mention gaze that can be approximated with head movement tracking, that another researcher in the group is investigating (see the instrumented caps in Figure 1).

Finally, as we had mentioned in the introduction, our goal was to evaluate our assistant in a realistic user study, and not jump into experiments with Large Language Models. However, we have started experiments in that respect, both as concerns the specific modules in our pipeline (for example, the embeddings of the semantic slots) and the system as a whole. So far, we have noticed that while ChatGPT (released exactly after we finished the COVID user study) is able to generate charts in response to specific language instructions, if appropriately connected to visualization software, it is not able to resolve referring expressions, i.e., to create a new visualization whose specification is partly derived from the referent. But this will be the topic of a future paper.

Acknowledgments

This work is supported by awards 2007257 and 2008986 from the National Science Foundation.

References

- Christopher Andrews, Alex Endert, Beth Yost, and Chris North. 2011. Information visualization on large, high-resolution displays: Issues, challenges, and opportunities. *Information Visualization*, 10(4):341–355.
- Jillian Aurisano, Abhinav Kumar, Alberto Gonzales, Khairi Reda, Jason Leigh, Barbara Di Eugenio, and Andrew Johnson. 2015. "Show me data": Observational study of a conversational interface in visual data exploration. In *IEEE VIS*, volume 15, pages 1–2.
- Jillian Aurisano, Abhinav Kumar, Alberto Gonzalez, Jason Leigh, Barbara Di Eugenio, and Andrew Johnson. 2016. Articulate2: Toward a conversational interface for visual data exploration. In *IEEE Visualization*, volume 8.
- Kenneth C. Cox, Rebecca E. Grinter, Stacie Hibino, Lalita Jategaonkar Jagadeesan, and David Mantilla. 2001. A multi-modal natural language interface to an information visualization environment. *International Journal of Speech Technology*, 4:297–314.
- Jacob Eisenstein and Randall Davis. 2006. Gesture improves coreference resolution. In *Proceedings of the Human Language Technology Conference of the NAACL, Companion Volume: Short Papers*.
- Tong Gao, Mira Dontcheva, Eytan Adar, Zhicheng Liu, and Karrie G. Karahalios. 2015a. [Datatone: Managing ambiguity in natural language interfaces for data visualization](#). In *Proceedings of the 28th Annual ACM Symposium on User Interface Software and Technology*, UIST '15, page 489–500, New York, NY, USA. Association for Computing Machinery.
- Tong Gao, Mira Dontcheva, Eytan Adar, Zhicheng Liu, and Karrie G Karahalios. 2015b. [Datatone: Managing ambiguity in natural language interfaces for data visualization](#). In *Proceedings of the 28th Annual ACM Symposium on User Interface Software & Technology*, pages 489–500. ACM.
- Enamul Hoque, Vidya Setlur, Melanie Tory, and Isaac Dykeman. 2017. Applying pragmatics principles for interaction with visual analytics. *IEEE transactions on visualization and computer graphics*, 24(1):309–318.
- Julie Hunter. 2014. Structured contexts and anaphoric dependencies. *Philosophical Studies*, 168:35–58.
- Ryu Iida, Masaaki Yasuhara, and Takenobu Tokunaga. 2011. Multi-modal reference resolution in situated dialogue by integrating linguistic and extra-linguistic clues. In *Proceedings of 5th International Joint Conference on Natural Language Processing*, pages 84–92.
- Qi Jiang, Guodao Sun, Yue Dong, and Ronghua Liang. 2021. Dt2vis: A focus+ context answer generation system to facilitate visual exploration of tabular data. *IEEE Computer Graphics and Applications*, 41(5):45–56.
- Hansol Kim, Kun Ha Suh, and Eui Chul Lee. 2017. Multi-modal user interface combining eye tracking and hand gesture recognition. *Journal on Multimodal User Interfaces*, 11(3):241–250.
- Michael Kipp. 2001. Anvil-a generic annotation tool for multimodal dialogue. In *Seventh European Conference on Speech Communication and Technology*.
- Michael Kipp. 2014. Anvil: The video annotation research tool. In Jacques Durand, Ulrike Gut, and Gjert Kristoffersen, editors, *Handbook of Corpus Phonology*, pages 420–436. Oxford University Press.
- Abhinav Kumar. 2022. *Towards a Context-Aware Intelligent Assistant for Multimodal Exploratory Visualization Dialogue*. Ph.D. thesis, University of Illinois Chicago.
- Abhinav Kumar, Jillian Aurisano, Barbara Di Eugenio, Andrew Johnson, Alberto Gonzalez, and Jason Leigh. 2016. [Towards a dialogue system that supports rich visualizations of data](#). In *Proceedings of the 17th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 304–309, Los Angeles. Association for Computational Linguistics.
- Abhinav Kumar, Jillian Aurisano, Barbara Di Eugenio, and Andrew E Johnson. 2020. Intelligent assistant for exploring data visualizations. In *FLAIRS Conference*, pages 538–543.
- Abhinav Kumar, Barbara Di Eugenio, Jillian Aurisano, Andrew Johnson, Abeer Alsaiani, Nigel Flowers, Alberto Gonzalez, and Jason Leigh. 2017. Towards multimodal coreference resolution for exploratory data visualization dialogue: Context-based annotation and gesture identification. In *The 21st Workshop on the Semantics and Pragmatics of Dialogue (SemDial 2017–SaarDial)(August 2017)*, volume 48.
- Xiaolong Li and Kristy Boyer. 2016. Reference resolution in situated dialogue with learned semantics. In *Proceedings of the 17th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 329–338.
- Lars Lischke, Lena Janietz, Anna Beham, Hartmut Bohnacker, Ulrich Schendzielorz, Albrecht Schmidt, and Paweł W Woźniak. 2020. Challenges in designing interfaces for large displays: the practitioners' point of view. In *Proceedings of the 11th Nordic Conference on Human-Computer Interaction: Shaping Experiences, Shaping Society*, pages 1–6.
- Arpit Narechania, Arjun Srinivasan, and John Stasko. 2020. N14dv: A toolkit for generating analytic specifications for data visualization from natural language queries. *IEEE Transactions on Visualization and Computer Graphics*, 27(2):369–379.

- Costanza Navarretta. 2011. Anaphora and gestures in multimodal communication. In *Proceedings of the 8th Discourse Anaphora and Anaphor Resolution Colloquium (DAARC 2011)*, Faro, Portugal, Edicoes Colibri, pages 171–181. Citeseer.
- Zahar Prasov and Joyce Y Chai. 2008. What’s in a gaze?: the role of eye-gaze in reference resolution in multimodal conversational interfaces. In *Proceedings of the 13th international conference on Intelligent user interfaces*, pages 20–29. ACM.
- Shaolin Qu and Joyce Y Chai. 2008. Beyond attention: the role of deictic gesture in intention recognition in multimodal conversational interfaces. In *Proceedings of the 13th international conference on Intelligent user interfaces*, pages 237–246. ACM.
- Norbert Reithinger, Dirk Fedeler, Ashwani Kumar, Christoph Lauer, Elsa Pecourt, and Laurent Romary. 2005. Miamm—a multimodal dialogue system using haptics. In *Advances in Natural Multimodal Dialogue Systems*, pages 307–332. Springer.
- Franca Rupperecht, Carol Naranjo, Achim Ebert, Joseph Olakumni, and Bernd Hamann. 2019. When bigger is simply better after all: Natural and multi-modal interaction with large displays using a smartwatch. In *Proceedings of the Twelfth International Conference on Advances in Computer-Human Interactions (ACHI 2019)*.
- David Schlangen, Timo Baumann, and Michaela Atterer. 2009. [Incremental reference resolution: The task, metrics for evaluation, and a Bayesian filtering model that is sensitive to disfluencies](#). In *Proceedings of the SIGDIAL 2009 Conference*, pages 30–37, London, UK. Association for Computational Linguistics.
- Ethan Selfridge, Iker Arizmendi, Peter Heeman, and Jason Williams. 2011. [Stability and accuracy in incremental speech recognition](#). In *Proceedings of the SIGDIAL 2011 Conference*, pages 110–119, Portland, Oregon. Association for Computational Linguistics.
- Vidya Setlur, Sarah E. Battersby, Melanie Tory, Rich Gossweiler, and Angel X. Chang. 2016a. [Eviza: A natural language interface for visual analysis](#). In *Proceedings of the 29th Annual Symposium on User Interface Software and Technology, UIST ’16*, page 365–377, New York, NY, USA. Association for Computing Machinery.
- Vidya Setlur, Sarah E Battersby, Melanie Tory, Rich Gossweiler, and Angel X Chang. 2016b. [Eviza: A natural language interface for visual analysis](#). In *Proceedings of the 29th Annual Symposium on User Interface Software and Technology*, pages 365–377. ACM.
- Leixian Shen, Enya Shen, Yuyu Luo, Xiacong Yang, Xuming Hu, Xiongshuai Zhang, Zhiwei Tai, and Jianmin Wang. 2023. [Towards natural language interfaces for data visualization: A survey](#). *IEEE Transactions on Visualization and Computer Graphics*, 29(6):3121–3144.
- Arthur Sluyters, Quentin Sellier, Jean Vanderdonckt, Vik Parthiban, and Pattie Maes. 2022. Consistent, continuous, and customizable mid-air gesture interaction for browsing multimedia objects on large displays. *International Journal of Human–Computer Interaction*, pages 1–32.
- Arjun Srinivasan and John Stasko. 2017. Orko: Facilitating multimodal interaction for visual exploration and analysis of networks. *IEEE transactions on visualization and computer graphics*, 24(1):511–521.
- Una Stojnic, Matthew Stone, and Ernie Lepore. 2013. [Deixis \(even without pointing\)](#). *Philosophical Perspectives*, 27(1):502–525.
- Yiwen Sun, Jason Leigh, Andrew Johnson, and Sangyoon Lee. 2010. Articulate: A semi-automated model for translating natural language queries into meaningful visualizations. In *International Symposium on Smart Graphics*, pages 184–195. Springer.
- Roderick S. Tabalba, Nurit Kirshenbaum, Jason Leigh, Abari Bhattacharya, Veronica Grosso, Barbara Di Eugenio, Andrew E. Johnson, and Moira Zellner. 2023. An investigation into an always listening interface to support data exploration. *Proceedings of the 28th International Conference on Intelligent User Interfaces*.
- Roderick S. Tabalba, Nurit Kirshenbaum, Jason Leigh, Abari Bhattacharya, Andrew E. Johnson, Veronica Grosso, Barbara Maria Di Eugenio, and Moira Zellner. 2022. [Articulate+ : An always-listening natural language interface for creating data visualizations](#). *Proceedings of the 4th Conference on Conversational User Interfaces*.
- Ryuichi Takanobu, Qi Zhu, Jinchao Li, Baolin Peng, Jianfeng Gao, and Minlie Huang. 2020. [Is your goal-oriented dialog model performing really well? empirical analysis of system-wise evaluation](#). In *Proceedings of the 21th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 297–310, 1st virtual meeting. Association for Computational Linguistics.
- Bonnie Lynn Webber and Breck Baldwin. 1992. [Accommodating context change](#). In *30th Annual Meeting of the Association for Computational Linguistics*, pages 96–103, Newark, Delaware, USA. Association for Computational Linguistics.
- Shomir Wilson, Alan W Black, and Jon Oberlander. 2016. This table is different: A wordnet-based approach to identifying references to document entities. In *Proceedings of the 8th Global WordNet Conference (GWC)*, pages 432–440.
- B. Yu and C. T. Silva. 2020. [Flowsense: A natural language interface for visual data exploration within a dataflow system](#). *IEEE Transactions on Visualization and Computer Graphics*, 26(01):1–11.