



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

Contact: abhatt62@uic.edu



Abari Bhattacharya¹, Abhinav Kumar¹, Barbara Di Eugenio¹, Roderick Tabalba², Jillian Aurisano³, Veronica Grosso¹, Andrew Johnson¹, Jason Leigh² and Moira Zellner⁴

- University of Illinois Chicago
- University of Hawaii at Manoa

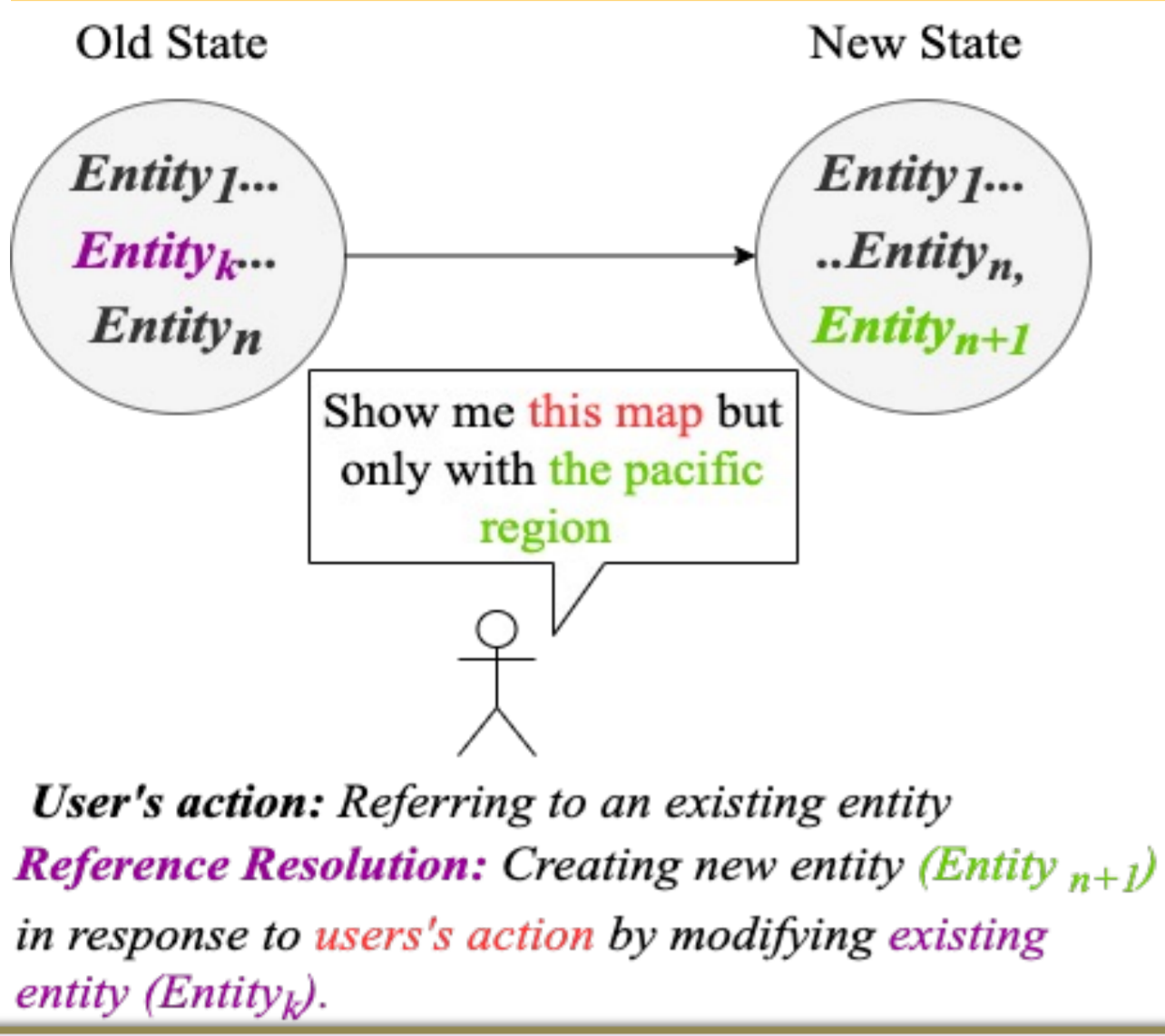


- University of Cincinnati
- Northeastern University



This work is funded by awards 200757 and 2008986 from the National Science Foundation

Introduction: Reference Resolution and New Entity Creation



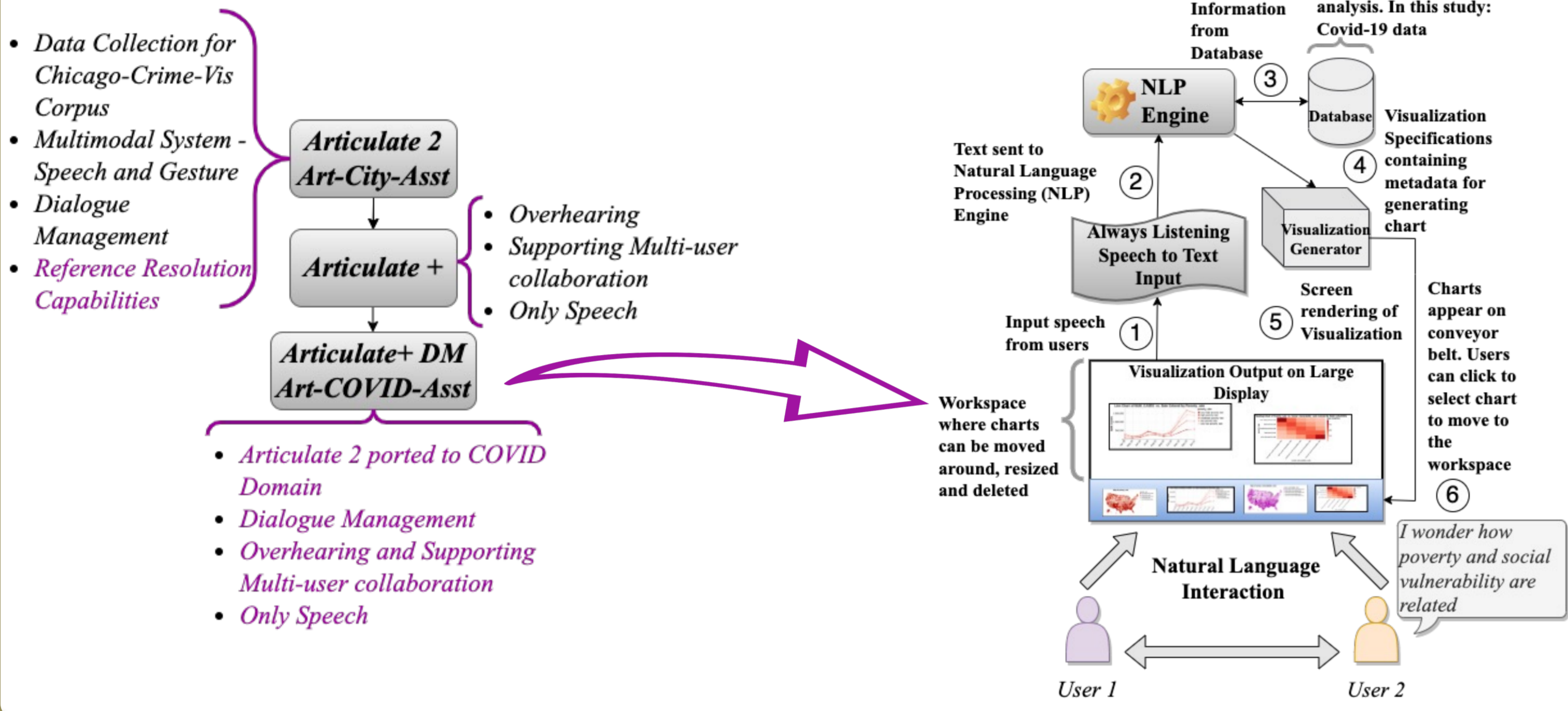
Our broader goal..

- Develop and deploy flexible conversational assistants to aid users explore data through visualizations
- Support user collaborations in exploratory data analysis

In this work we..

- Focus on reference resolution and new entity establishment
- Discuss reference detection and resolution in
 - controlled offline setting and
 - challenges presented when system is deployed "in the wild"

The "Articulate" Project: Background and Architecture



User Study with Art-COVID-Asst



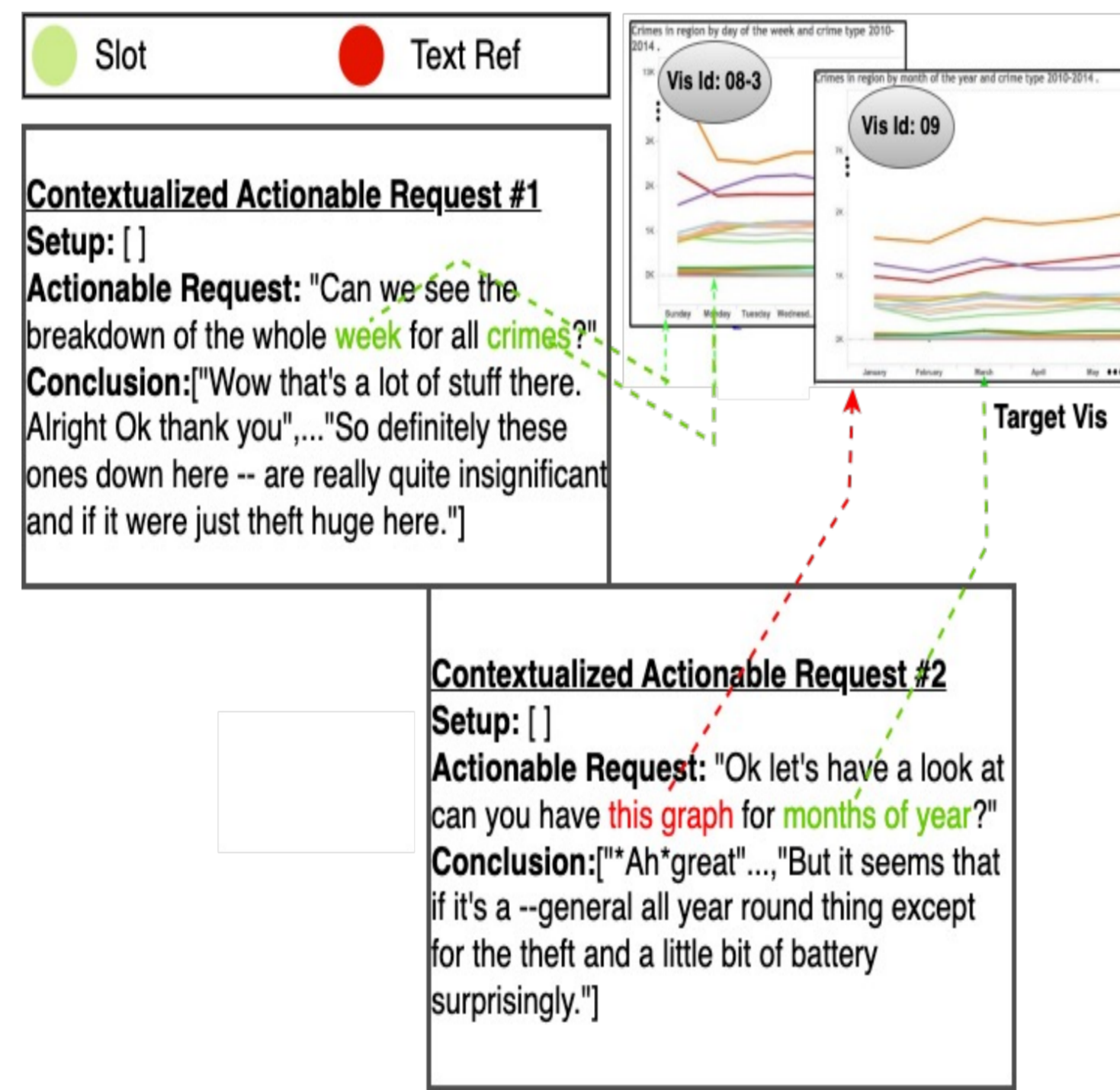
Porting to new domain: Art-COVID-Asst	<ul style="list-style-type: none"> Updating KO DM, Reference Detection, Reference Resolution: same as Art-City-Asst
User Study	<ul style="list-style-type: none"> 15 groups of 2 2 open-ended timed exploratory data analysis tasks pertaining to COVID mitigation
COVID Knowledge Ontology	<ul style="list-style-type: none"> 710 terms categorized into 13 semantic slots pertaining to COVID

Controlled vs Unconstrained: Evaluation and Results

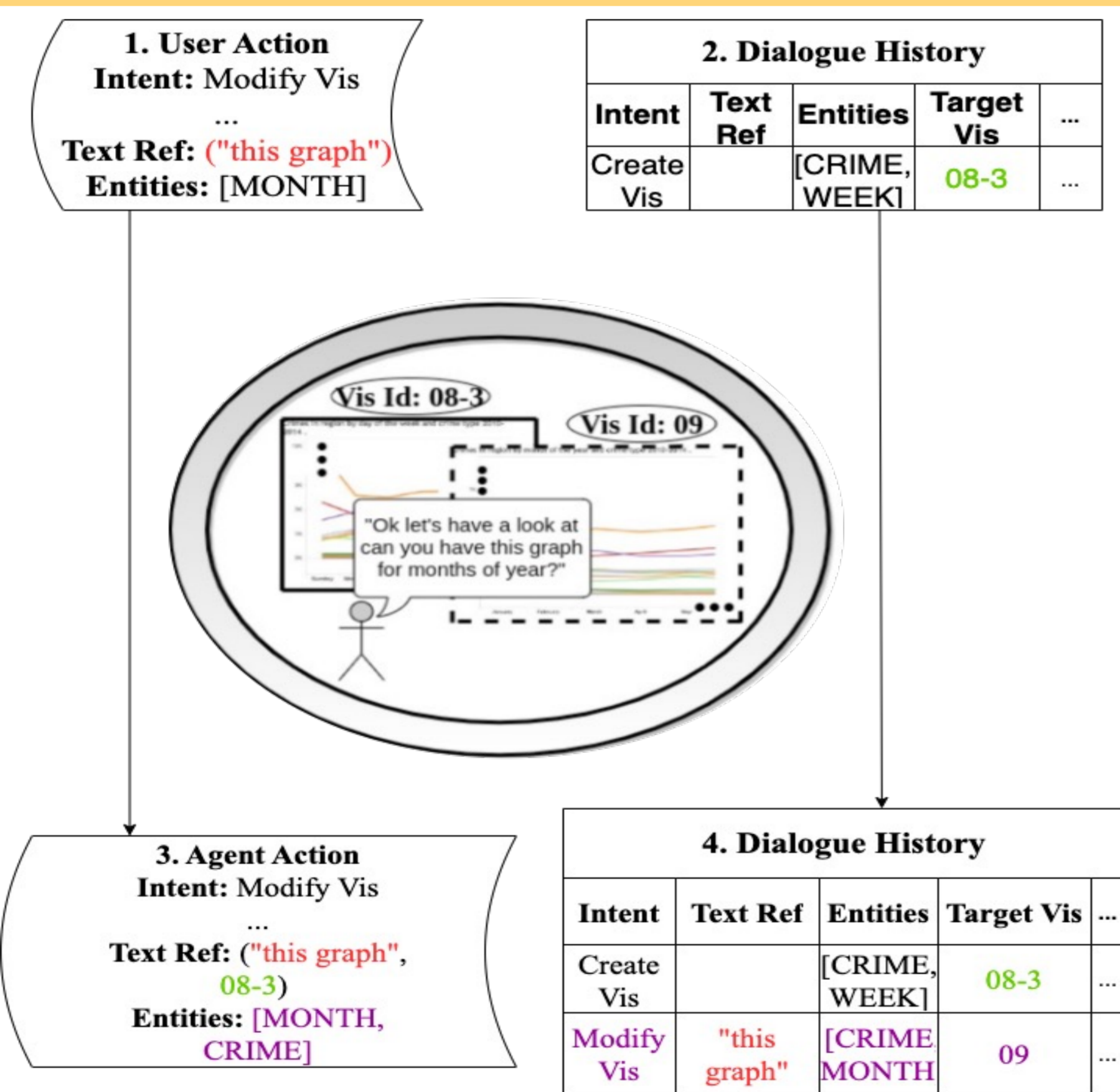
Evaluation Methodology	Controlled Setup Art-City-Asst	Unconstrained Setup Art-Covid-Asst	Remarks	
	<ul style="list-style-type: none"> Offline evaluation run on transcripts of Chicago-Crime-Vis corpus Focus on references occurring in setup and AR (specifically DA-s CREATEVIS and MODIFYVIS) Focus on single referents and single targets 	<ul style="list-style-type: none"> Real-time evaluation Speech recognition errors were a major bottleneck in the user study Experiments with transcripts of the user studies generated using Whisper speech recognition model and fed to the back-end code → COVID(T) – Transcripts (#utterances: 3096) Real time logs → COVID(A)- Automatic (# utterances: 8440) 	<ul style="list-style-type: none"> For evaluation of Art-COVID-Asst we need to manually verify the results returned by the reference pipeline. A significant sample size is computed for both Random significant sample of 340 (11%) utterances for COVID (A), and of 370 (4.38%) utterances for COVID (T). 	
		COVID(A)	COVID(T)	
Semantic Frame Accuracy	% of Visualization Frames			
	Correctly identified slots			
	All	55%	54%	52%
	At least 75%	85%	60%	63%
	None	7%	18%	16%
Reference Detection	Accuracy			
	Actionable Request	55.0%	25.0%	45.8%
Reference Resolution	Accuracy for Window sizes			
	1	74.4%	-	36.3%
	∞	68.3%	-	54.0%

Corpus Creation: Chicago-Crime-Vis

User Study	<ul style="list-style-type: none"> 16 participants interacting with human Visualization Expert (VE) 3.2K Utterances
Utterance Types	<p>Contextual Actionable Requests (CAR)</p> <ul style="list-style-type: none"> Setup: Think aloud prior to an actionable request for VE AR: The actionable request Conclusion: Think aloud after AR. 449 CARs covering 1545 utterances AR annotated with Dialogue Acts (DA): Notably, CREATEVIS (creating new visualization from scratch), MODIFYVIS (creating new visualization based on existing visualization)
Referring Expression Annotation	<p>294 References: 176 Textual</p> <p>680 Phrases as slot fillers corresponding to data attributes in the Knowledge Ontology</p> <p>Inter-annotator Agreement: $\kappa = 0.85$</p>
Chicago-Crime-Vis Knowledge Ontology	<ul style="list-style-type: none"> Constructed semi-automatically 3.5K terms categorized into 11 semantic pertaining to the city of Chicago



Co-reference: Detection, Resolution and New Entity Establishment



Semantic Frame Construction (CREATEVIS and MODIFYVIS)	<ul style="list-style-type: none"> "month of year" → MONTH in (1)
Reference Detection	<ul style="list-style-type: none"> Sequence Tagging (IOB2 format) CRF Model with POS tags as features (F1 = 61.2% on B-REF, I-REF, O-REF task) "this graph" in (1)
Reference Resolution (Heuristics based on Similarity)	<ul style="list-style-type: none"> Candidate visualization with highest cosine similarity (threshold of 0.4, empirically established) selected. "08-3" in (2)
New Entity Establishment (New visualization)	<ul style="list-style-type: none"> new visualization is constructed ("09") with referent's frame representation to infer missing information ("08-3") → Agent action (3) updated Dialogue History (4)

Discussion

Findings from Evaluation

Semantic Frame Accuracy	<ul style="list-style-type: none"> We report partial accuracy to provide more nuanced analysis of the assistant's performance In dialogue-based application for data exploration like ours, partially recognized VF can generate charts This may help the users move forward. Irrespective of subpar performance of speech to text algorithm more than 60% VFs had 75% or more slots correctly filled Attested by questionnaires filled by users post user study <ul style="list-style-type: none"> Mean scores of 4 and 3 respectively for usefulness of generated charts and ease of command system use on a 5-point Likert scale
Reference Detection	<ul style="list-style-type: none"> Unlike controlled study setting with one subject, when two people collaborate for exploratory task, three things happen. <ol style="list-style-type: none"> They talk to each other Make requests to the system Finally draw conclusions. Reference detection in real-time utterances extremely complex. In the case of COVID (A), we also attribute the lack of accuracy to speech-recognition errors.
Reference Resolution	<ul style="list-style-type: none"> Speech recognition error major roadblock for lack of resolved references in COVID(A) In the unconstrained setting, we observe when two people are involved in the conversation, there are more relevant entries in DH → expanding the window yields better results

Future Work

- Modeling user behavior for referring more distant visualization
- Leverage multi-modality – gesture, eye gaze and head movement tracking, etc.
- Experiments with Large Language Models