

Recommendations for Visualization Recommendations: Exploring Preferences and Priorities in Public Health

Calvin Bao
csbao@umd.edu
University of Maryland
College Park, MD, United States

Siyao Li
siyaoli@terpmail.umd.edu
University of Maryland
College Park, MD, USA

Sarah Flores
sflores3@terpmail.umd.edu
University of Maryland
College Park, MD, USA

Michael Correll
mcorrell@tableau.com
Tableau Research
Seattle, WA, USA

Leilani Battle
leibatt@cs.washington.edu
University of Washington
Seattle, WA, USA

ABSTRACT

The promise of visualization recommendation systems is that analysts will be automatically provided with relevant and high-quality visualizations that will reduce the work of manual exploration or chart creation. However, little research to date has focused on what analysts *value* in the design of visualization recommendations. We interviewed 18 analysts in the public health sector and explored how they made sense of a popular in-domain dataset¹ in service of generating visualizations to recommend to others. We also explored how they interacted with a corpus of both automatically- and manually-generated visualization recommendations, with the goal of uncovering how the design values of these analysts are reflected in current visualization recommendation systems. We find that analysts champion simple charts with clear takeaways that are nonetheless connected with existing semantic information or domain hypotheses. We conclude by recommending that visualization recommendation designers explore ways of integrating context and expectation into their systems.

CCS CONCEPTS

• **Human-centered computing** → **Visualization design and evaluation methods; Visualization systems and tools; User studies.**

KEYWORDS

Visualization recommendation systems, algorithmic trust, automation, recommendation source

ACM Reference Format:

Calvin Bao, Siyao Li, Sarah Flores, Michael Correll, and Leilani Battle. 2022. Recommendations for Visualization Recommendations: Exploring Preferences and Priorities in Public Health. In *CHI Conference on Human Factors*

¹National Health and Nutrition Examination Study 2013-2014 [3].

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

CHI '22, April 29-May 5, 2022, New Orleans, LA, USA

© 2022 Copyright held by the owner/author(s). Publication rights licensed to ACM.
ACM ISBN 978-1-4503-9157-3/22/04...\$15.00
<https://doi.org/10.1145/3491102.3501891>

in *Computing Systems (CHI '22)*, April 29-May 5, 2022, New Orleans, LA, USA.
ACM, New York, NY, USA, 17 pages. <https://doi.org/10.1145/3491102.3501891>

1 INTRODUCTION

Data analytics, and especially the creation of informative and useful visualizations of large datasets, can be a time-consuming and complex process. As part of a larger design goal of “augmenting” analytics to offload effort to algorithmic systems [14], there is a growing number of systems that automatically generate and recommend visualizations [49, 52]. Systems such as Voder [40], Dziban [29], and Tableau Show Me [31] can generate visualizations to both surface potentially insightful features of a dataset as well as provide guidance for “novice investigators” to generate their own visualizations [22].

The designers of recommendation systems have explicit or implicit **design values** about what charts they surface: for instance, recommenders that purport to automatically surface “insights” [4] might place value on particular statistical patterns like outlying values or correlated fields [9, 32]. However, despite the proliferation of visualization systems in the literature [49, 52], there has been little work on interrogating these design values, and observing matches and mismatches between the values of recommendation system *designers* and *consumers*. Our work is therefore focused on a central question: **what sort of visualizations do people want to see, and how well do these preferences actually align with the sorts of visualizations that algorithmic recommendation systems currently provide?**

Prior work considers how people react to different recommendation sources [36, 48], but does not consider the *priorities* and *expectations* analysts have when creating their own recommendations for other analysts. Without a deeper understanding of how analysts themselves think about the visualization recommendation process, new recommendation engines may barely help [49], and possibly even hinder [5], an analyst’s ability to explore their data, creating a “double-edged sword” [28] of potentially “opaque, inflexible, brittle, and domineering” [32] analysis.

In this paper, we present the results of a pre-registered² qualitative study designed to interrogate and elicit design values around generating and evaluating visualization recommendations. Our study, conducted with public health researchers supplied with a

²https://aspredicted.org/AEI_GBA

sample dataset of U.S. self-reported health data, consists of two components:

- (1) an **ideation task** where participants, with the help of an experienced visualization designer working with Tableau, sketched and then realized their own visualization recommendations for an imagined client seeking to influence public policy, and
- (2) a **selection and ranking task** where participants explored a gallery of recommendations (some generated automatically by systems, and some by human curators) and selected the ones they felt were most valuable for their client.

We chose these tasks and this participant pool to examine points of friction between the values of recommendation systems and analysts. I.e., we wanted to contrast the (often) domain-agnostic assumptions of visualization recommendation systems with the specific domain expertise and context of our participants, and contrast the (frequent) focus on narrowly defined statistical findings in recommendation systems with the unconstrained and diverse rhetorical and persuasive goals of our participants.

Of the design values we encountered in our exploration, the three most prominent that our participants valued in recommended visualizations were:

- (1) **simplicity**— participants, often with an assumed audience in mind, valued simple visualization designs over more complex ones, and visualizations with one clear takeaway over more nuanced or complex data stories. Titles and labels, filtering, and aggregation were common strategies to reduce the complexity of data.
- (2) **relevance**— in addition to a preference for the removal of extraneous data from recommendations, participants also made efforts to tailor their charts to their domain of interest. E.g., a preference for bivariate visualizations with anticipated casual relationships (e.g. that one variable would “drive” another, or produce a clear “trend”).
- (3) **interestingness**— participants were reluctant to provide visualizations that failed to show clear trends, group differences, or other strong signals. Participants wanted recommended charts to provide direct evidence for or against particular hypotheses, or to promote specific follow-up actions.

These design values suggest both *opportunities* and *dangers* for designers of future visualization systems. On the one hand, they suggest benefits for incorporating additional data semantics or explicit user intent into recommendation systems to better meet the goals of analysts. On the other hand, they suggest that care should be taken to communicate complex or ambiguous trends in the data that might arise in the recommendation process, and that the desire to surface strong signals promotes a form of exploratory data analysis that lends itself to false positives or other dangers [37, 50] to the reliability or robustness of findings.

2 RELATED WORK

Our research questions and experimental design are informed by assumptions and goals behind the design of existing visualization recommendation systems, as well as by prior studies that involve participants expressing their preferences amongst visualizations

from heterogeneous sources or creating novel and heterogeneous visualizations themselves. We therefore highlight three topics of related research: visualization recommendation systems, assessment of those systems (and visualizations in general), and visualization construction for novice users.

2.1 Visualization Recommendations

Visualization recommendation systems aim to ease the process of visualization authoring or exploratory data analysis for different user groups [17, 52]. Each system has its own set of metrics and structures to represent what users find valuable to visualize in a dataset [49]. For example, some recommendation systems prioritize perceptually effective encoding channels for a given set of data attributes (e.g., [29–31]), popular visualization designs that other users have created in the past (e.g., [15]), or specific types of data trends such as pairwise correlations between attributes [9] or significant differences among sub-populations in the dataset [43]. However, these priorities are often set in a way that is agnostic to either the domain of interest or the particular analytical goals of the user. When the system’s and user’s priorities are misaligned, the system may generate distracting and ineffective recommendations [49]. In this paper, we seek to clarify what analysts prioritize when designing their own recommendations in the context of public health, and to understand how analysts’ priorities compare with those of existing systems.

We divide existing visualization recommendation systems into three categories: Auto-Insight, Encoding, and Q&A, although we note that these categories are not necessarily mutually exclusive, and that recommendation systems can and do incorporate design values or patterns from multiple modalities.

2.1.1 Auto-Insight. Auto-insight systems automatically detect and visualize meaningful attributes, trends, or other statistical properties within a provided dataset [28], removing some of the labor or luck involved in manual exploratory data analysis [14]. These data insights can be given in the form of text describing statistical patterns in the data or through visualizations [25]. Example systems include Voder [40], which focuses on textual facts and insights, as well as PowerBI Quick Insights [34], Foresight [9], and Amazon QuickSight [1], which focus on insights presented as visualizations. Voder creates textual “data facts” based on the dataset’s attributes to assist users in interpreting generated data visualizations and communicating findings [40]. PowerBI’s Quick Insight [34] panel searches through different subsets of a dataset and detects particular classes of statistical features (e.g., outliers, variance, correlation, categories with a strong majority) to generate insights during data exploration. Similarly, Foresight [9] ranks visualizations based on statistical properties present in the data. Lastly, Amazon QuickSight [1], within the broader Amazon Web Services ecosystem, also creates data summaries using in-house algorithms, allowing users to upload and integrate their own models.

Although each of these systems provides a different means of exploring data and communicating insights, they generally lack explicit explanations for how insights are generated, leading some users to distrust the results [49]. Furthermore, users do not necessarily know whether these recommendations cover everything that could or should be learned from the given dataset [25]. This

lack of transparency and consideration of user context (i.e., user preferences and intended recommendation goals) may result in bias, unreliability, and disruption of the exploratory data analysis process [28, 48]. Another worry is that, by exhaustively searching for potentially interesting patterns, auto-insight systems can function as “p-hacking machines” [5, 37], surfacing “insights” that are ultimately spurious or misleading.

Our experimental design is most closely aligned with the goals and values of auto-insight recommenders, in that our participants were asked to generate meaningful visualizations for their clients without constraints on fields or designs of interest, although we note overlaps with other forms of recommenders below.

2.1.2 Encoding. We define encoding recommendation systems as those that suggest designs of individual visualizations given user-specified data attributes. These recommendation(s) are often based on the characteristics of the data and expert knowledge on the expressiveness and effectiveness of different encoding channels or chart designs [30].

These systems employ a variety of approaches in how they encode expert knowledge. Draco [35] uses a set of constraints to assist users in visualization design and prioritize visual exploration, promoting effective encodings, and predicting the best visualization through a ranking system. Dziban [29] builds upon the Draco knowledge base while incorporating chart similarity logic to create a balance between “automated suggestions and user intent.” Show Me [31] either suggests (and automatically generates) particular chart designs given the data types of selected data attributes, or allows the user to progressively construct a chart by adding attributes one at a time, automatically suggesting new encodings or chart designs. Graphscape [19] uses a directed graph model in which nodes represent chart specifications and edges represent transitions between charts. This model enables Graphscape to recommend alternative designs by minimizing the perceptual distance between new recommendations and visualizations previously seen. Table2Charts [51] takes table-chart pairs and learns patterns that assist in generating recommendations. Several systems also focus on multi-dataset exploration. For example, GEViTRec [7] recommends visualizations across multiple datasets by looking for linking fields or domain-centric constraints. Lastly, Data2Vis [10] uses neural networks to translate a given dataset into a resulting visualization specification, based on a training set of presumably well-designed Vega-Lite [38] specifications. VizML [15] applies a similar learning approach to Plotly visualizations.

While our experimental framing was less aligned with the design values of these systems, as our participants had free choice over which variables to include, the ability of our participants to create and select their own designs, and to iteratively alter the default designs generated over the course of the experiment, allowed us to see if existing assumptions around expressiveness and effectiveness matched the preferences and priorities of our participants (who, while embedded in their domain of interest, had varying levels of expertise in visualization design).

2.1.3 Q&A systems. We consider Q&A recommenders to be systems where the recommendation engine and the user can engage in one or more rounds of communication for the generation and refinement of recommendations. A prototypical Q&A system might take

as input textual questions, suggestions, and/or attributes (“questions”) and produce as output an appropriate visualization (an “answer”). Amazon QuickSight and Tableau Ask Data are examples of systems that provide this feature. QuickSight [1], for instance, features a search bar wherein users can enter natural language questions about their data. The user’s intent is then inferred from the questions, and the system returns an answer in the form of a number, visualization, or table. Tableau’s Ask Data [46], a system for performing ad-hoc exploration and analysis, also incorporates natural language interaction features: users type in natural statements or questions into an input bar and the system produces a chart [41].

Although the user can explicitly tell these Q&A systems what they are interested in, the extent to which systems are truly cognizant of or reactive to the intent of the user is often unclear [41]. For example, a user may have a chart in mind when asking a question (or selecting attributes) but receive an entirely different chart as output, deviating from their expectations. How should systems respond to ambiguous questions from the user? [13]

In our study, we are interested in understanding what users generally value in the design and construction of visualizations, which can inform general purpose guidelines for creating intent-focused Q&A systems, and other visualization recommendation systems. Our study protocol also allowed participants to iterate with us to refine their recommendations, affording an analysis of what sorts of refinement or repair operations are common in visualization recommendation with human partners, that could similarly be of use for designers of automated Q&A systems.

2.2 Visualization Recommendation Assessment

While a full consideration of all of the ways visualizations can and have been assessed is out of the scope of this work (see Lam et al. [23] for a typology), we focus on studies dealing with eliciting preferences from sets of unfamiliar visualizations that have been presented to participants, especially in the context of *recommendation*. Peck et al. [36] performed a qualitative study where participants were asked to assess their attitudes towards an array of infographics, highlighting how different beliefs and stances can influence the perceived quality, utility, and trustworthiness of a visualization. Lee et al. [27] explored the process of making sense of unfamiliar visualizations through a think-aloud procedure similar to the one we adopt, with an emphasis on investigating what factors influence the interpretability of a visualization. In our study, we look to these works as a model for exploring attitudes towards *existing* visualizations, but we also include a visualization *authoring* step in order to assess specific design characteristics our participants valued when creating their own recommendations.

Another paper we use as a model is Zehring and Singhal et al. [48], where participants were given sets of recommendations from an unseen visualization recommender, with the specific goal of evaluating how the stated provenance of the recommender (human or algorithmic) impacted perceived quality and trust. While we similarly intermix human and algorithmic recommendations in our study with the goal of investigating any systematic differences between the two sources, our work moves away from questions

of trust and provenance and towards broader issues of perceived utility and impact. Lastly, Zeng et al. [49] propose a framework for specifying multiple visualization recommendation algorithms within the same semantic space to enable quantitative comparison and evaluation. While our mixture of qualitative and quantitative methods complicate this process, our study findings could be incorporated into the Zeng et al. framework to improve evaluation of end-user preferences and expectations for recommended visualizations.

2.3 Barriers and Methods for Eliciting Visualizations

As we intended for our participants to both evaluate existing visualization recommendations *and* generate their own, we explored potential processes and pitfalls for eliciting visualizations from diverse audiences, especially audiences who may lack experience with existing visualization design tools. Several visualization construction barriers exist for visualization users (especially novices) – Grammel et al. [11] finds that novices struggle with navigating and mapping the relationships between visualization concepts: exploratory questions, data attributes, and visualizations during the construction process. Other barriers include several reported by Kwon et al. [22]: a failure to interpret visualizations properly and a failure to match expectations and functionality of the visualization. These barriers often caused frustrations among novice visualizers. This presents an underlying “gulf of execution” between the types of visualizations that users want versus what visualization recommendation systems actually generate. For this work, we are particularly interested in how more advanced analytics users approach this gulf.

To help reduce the barriers to effective visualization construction, free-form sketching can serve as an expressive medium of converting internal thought to external representations [21, 45]. Moreover, work by Tversky highlights the power of the sketching process to reveal the designer’s underlying ideas and reflect core aspects of one’s prioritization [42], a power that is used by systems such as SketchStory [26] for fluid and flexible visualization authoring. A tangential effect of free-form sketching is that it provides direct interaction. Studies on whiteboard usage showed how whiteboard sketching enables people to immediately externalize ideas without being interrupted by or having to translate their ideas to another medium or system [44]. We incorporated free-form sketching into our study to reduce the construction complexity for our participants and to better observe expectations for the visualizations they create.

2.4 Summary

Our analysis of prior work points to a wide space of visualization recommendation systems that nevertheless prioritize specific statistical features, low-level analysis tasks, and visualization design rules, all of which have advantages in particular scenarios, but may or may not capture the specific priorities and mental models of analysts more broadly. Prior work also suggests relevant strategies for working with audiences across levels of data expertise or engagement to develop rich frameworks around understanding, values, and priorities in visualization. Our study seeks to integrate these two perspectives by performing a human-centric assessment

of the priorities and values of visualization recommendations, the results of which can guide the designers of future visualization recommendation systems.

3 MOTIVATION

Our study is motivated by a potential gap in *design values*: between the values of designers of visualization recommendation systems (who might prioritize highlighting a particular subset of statistical patterns, data facts, or “insights”) and those of human analysts (who might have more semantically rich or teleological expectations of their visualizations). With a deeper understanding of what analysts prioritize as they create and rank visualizations for later recommendation, we can compare our observations with how visualization recommendation systems are currently designed, and provide concrete feedback for how current and future systems can be refined to more closely align with the goals and values of their end-users.

We break our broader research question (*what do analysts value in the design of visualization recommendations, and are these design values reflected in current visualization recommendation systems?*) down into three sub-questions to investigate through our study:

- **RQ1:** What characteristics of a visualization design do analysts prioritize when recommending them to colleagues?
- **RQ2:** What do analysts prioritize when evaluating visualization recommendations from other sources?
- **RQ3:** How do the recommendations made by analysts align with those created from other sources in terms of visual form or analytical purpose?

While recommendation systems are often agnostic or insensitive to data domain or analytic intent, our belief is that the perceived usefulness of a visualization is often task- and domain-dependent [39]. For these reasons, we focus on a single domain in this work in order to specifically elicit any potential tensions between the domain insensitivity of many automatic recommenders and the domain knowledge and intents of our participant pool. Specifically, we investigate how researchers and professional analysts working in the public health sector create and evaluate visualization recommendations for a goal of presenting information to shape public policy.

Our questions are ones of exploring or enumerating *alignment* in design values rather than evaluating predictions or building models. As such, we do not enumerate hypotheses for testing, but focus more on descriptive quantitative reports of our findings augmented with qualitative data.

4 EXPERIMENT DESIGN

We designed a pre-registered³ experiment to better understand what visual or data characteristics analysts prioritize when creating visualizations for other analysts, and what analysts purport to value when presented with a gallery of human-curated and algorithmically-generated recommendations.

To accommodate a wider range of participants as well as to abide by COVID-19 pandemic protocols, the study was conducted online using the video conferencing platform Zoom. Participants shared their screen with the experimenters, and completed the study using Google Jamboard, an online sketching and whiteboard tool. In the

³<https://aspredicted.org/7d7gd.pdf>

following subsections, we describe our participant pool, pilot study, and visualization artifacts used for the experiment, and then walk through the entirety of an interview, describing each phase and how the participant was to interact with the interviewing team.

Additional study details, including transcripts, sketches, generated analyses, data tables, and analyses are available at <https://osf.io/xeub3/>.

4.1 Participants

After approval by our institutional IRB, we recruited 18 participants through a combination of university mailing lists, snowball sampling through research collaborators, and advertising on social media. We employed different methods of sampling to broaden population groups of participants to minimize the selection bias in our recruitment process.

We present demographic information about our participants in Table 1. Our participants ranged between 18-64 years old, with two being between 19-24, ten being between 25-34 years old, three being 35-44 years old, and three being 45+ years old. In terms of domain expertise, at the time of study, four participants were current graduate students and the remaining fourteen participants were working as public health professionals in various capacities, ranging in roles from project director, health program administrator, faculty member, and research scientist. Regarding frequency of creating data visualizations, nine participants reported creating visualizations at least once a month, five reported creating visualizations at least once a week, and three reported creating visualizations daily. The remaining participant reported creating visualizations rarely (less than once a month). To qualify for participation, participants had to have at least two years of industry analyst or research experience in public health. We compensated participants with a \$25 Amazon gift card for completing the study.

By the end of the study, we collected a set of 53 visualization sketches (all participants but one sketched out three, while the one sketched only two) and 18 rankings of the visualizations in a gallery of recommendations from our participants.

4.2 Experimental Dataset

To ensure that we selected and presented data that aligned with the interests of our target participants, we solicited feedback from experts in public health at our primary authors' home institution. These experts provided guidance on relevant datasets, attributes that would be of particular interest to a public health audience, and which groups and departments to target for recruitment. Based on this feedback, we selected a vertical subset of the National Health and Nutrition Examination Survey (NHANES) from 2013-2014 for use in our study [3]. Twenty attributes were extracted from the NHANES dataset, covering the *Demographics*, *Examinations*, *Dietary*, and *Questionnaire* response categories. We randomized the ordering of attributes for each participant to mitigate order effects. A sample table of 6 records was provided to participants, so they could see the available attributes and their data types. Participants were also given the option to view the dataset in its entirety through an online link to a spreadsheet.

4.3 Pilot Study

We conducted an initial pilot experiment with five participants. We also presented our experimental protocol to two faculty members at our institution's School of Public Health for additional feedback.

Initially, we asked pilot participants to hypothesize about which attributes would be important for visualization. This helped us to narrow down our list of data attributes that we believed to be valuable and relevant to in-domain analysts. We specifically extracted only these data attributes from the broader dataset. We presented this subset to each pilot participant and asked them to sketch five visualizations they would recommend to other analysts that would explore the same dataset. We found that asking for five sketches prohibitively extended the length of the study, as they reported that it was difficult to create five sufficiently distinct and interesting visualizations in the allotted time. We decreased the number to three, resulting in a final approximate study length of 60 minutes. We validated this updated study design with an additional pilot participant. An additional modification as a result of this piloting was to allow both solicitation of sketches via Jamboard (which some participants found limiting or difficult to use) as well as via hand-drawn sketches emailed directly to the experimenters.

We also used the pilot study to seed the gallery of visualizations that we ultimately used for our selection/ranking task in the final experiment. Our gallery was created via a mixture of the visualizations created by our pilot participants as a result of their ideation task (5 visualizations) and the visualizations favored by our pilot participants from a gallery of visualizations generated by recommendations systems Dzuban [29], Voyager [47], PowerBI [34], Tableau ShowMe [31], and Amazon QuickSight [1]. We selected five of the algorithmically-recommended visualizations with the most votes from our pilot participants to add to our finalized gallery of recommendations: the resulting favorite visualizations came from Dzuban (3 visualizations) and Voyager (2 visualizations).

4.4 Experiment Flow

Our study consisted of four phases:

- **Phase 1. Tutorial** The participant receives a brief tutorial on how we use Jamboard to conduct the virtual experiment. The participant was then introduced to the task for Phases 2 and 3, and given an opportunity to review the dataset.
- **Phase 2. Visualization Ideation with Mediator and Wizard** The participant sketches 3 visualizations based on the given dataset and task in conversation with a mediator. A "wizard" then translates these sketches into final visualizations using Tableau Desktop.
- **Phase 3. Visualization Selection and Ranking** The participant selects and ranks 5 visualizations from a set of 10 visualizations (5 human-generated and 5 system-generated).
- **Phase 4. Post-survey** Participants complete a short survey to provide closing comments and demographic information.

4.4.1 Tutorial. We explained the features available on Jamboard (e.g., pencil, eraser, drawing shapes) Participants spent five minutes familiarizing themselves with Jamboard's features. They then were given the task, which is as follows:

"You have been paired with an analyst to develop slides to present to a client. The client is developing public policy to improve health

Table 1: Demographic information about each study participant, labeled by participant ID (PID).

PID	Age Group	Gender	Job Title	Perform Data Analysis	Create Visualizations
1	19 - 24	Female	Consultant	Weekly	At least once/month
2	25 - 34	Female	Research Coordinator	Weekly	At least once/week
3	25 - 34	Female	Graduate Student	<Once/month	At least once/month
4	25 - 34	Female	Project Coordinator	<Once/month	At least once/month
5	35 - 44	Female	Consultant	Weekly	At least once/week
6	25 - 34	Male	Research Coordinator	<Once/week	At least once/week
8	19 - 24	Female	Health Program Administrator	<Once/month	At least once/month
9	45 - 54	Female	Assistant Professor	Weekly	At least once/month
11	25 - 34	Female	Graduate Student	<Once/week	At least once/month
12	25 - 34	Female	Research Scientist	Daily	Daily
13	55 - 64	Female	Lecturer	<Once/week	At least once/month
14	35 - 44	Female	Faculty	Weekly	At least once/month
15	45 - 54	Female	Project Director	<Once/week	At least once/month
16	25 - 34	Female	Research Scientist	Daily	Daily
17	35 - 44	Male	Project Coordinator	<Once/week	At least once/week
18	25 - 34	Female	Faculty Specialist	Weekly	At least once/week
19	25 - 34	Female	Graduate Student	<Once/month	At least once/month
20	25 - 34	Male	Database Administrator	<Once/month	Daily

outcomes in the US. Create a set of 3 visualizations (a.k.a. charts and graphs) that you would recommend to this client.”

They were then given an opportunity to explore and review the dataset, with minimal input from the interviewers.

4.4.2 Visualization Ideation with Mediator and Wizard. Participants were then asked to select and rank a short list of attributes of interest to them, to better ground them within the task. They then were asked to sketch visualizations while following a think-aloud manner. The participant shared with the mediator their sketches and any further necessary visualization specifications, while discussing the ideas, goals, and motivations behind the visualization. Then, the wizard re-created the participants’ sketch with visualization software, Tableau Desktop, similarly to the process used by Grammel et al. [12] in their study. During this sub-phase, the mediator guided the participant by explaining set tasks, and organizing the jamboard. Concurrently, the wizard would share their Tableau visualizations with the participant, receiving feedback and then updating the visualizations based on this feedback to ensure that the final designs matches the participant’s sketched intentions.

4.4.3 Visualization Selection and Ranking. Afterward, participants were presented with a gallery of 10 visualizations (Figure 1, a mixture of human-recommended and algorithmically recommended visualizations (see subsection 4.3)). We did not specify the provenance of these visualizations. The display order of visualizations was randomized for each participant to account for order effects. Participants were asked by the mediator to select and rank the top 5 visualizations within this set that they favored and would recommend to others in relation to the prompt given in 4.4.2.

4.4.4 Post-survey. The last step of the experiment was for participants to fill out a short post-interview survey about their visualization experience, data analysis experience, age, and work experience,

e.g., their highest level of education, years of experience, prior statistical experience, current job title, and also the frequency for performing data analysis.

4.4.5 Data Collection. Participants’ survey responses were collected using online forms and stored as CSV files. Participants’ sketches (uploaded images and Jamboard designs), gallery selections and survey responses were also collected. The sketches were translated into corresponding Vega-Lite specifications [38]. The gallery selections and rankings were also stored as CSV file. Participants shared their screen during the experiment and screen capture was recorded, providing video and audio recording of each interview session. Each session was transcribed. Experimenters also took notes during the experiment. To preserve participant privacy we do not share these audio or video records directly, but only the transcripts, notes, and codes, and after a manual process of redacting identifying information. We share this collected data as a public resource on OSF: <https://osf.io/xeub3/>.

4.5 Experiment Design Limitations & Trade-offs

We considered multiple trade-offs in the design of our study, which we discuss here.

4.5.1 Limitations in Data Collection. In our experiment, we chose to reduce the number of attributes presented to participants for two reasons related to participant accommodation. First, we found in our pilot study that having access to the entire NHANES dataset was overwhelming for participants, and they reported encountering “analysis paralysis” when deciding what attributes to select and what visualizations to sketch. Second, with the original number of attributes to choose from, thousands of attribute combinations were possible for creating visualizations, which could have led to participants having little or no overlap in their attribute and

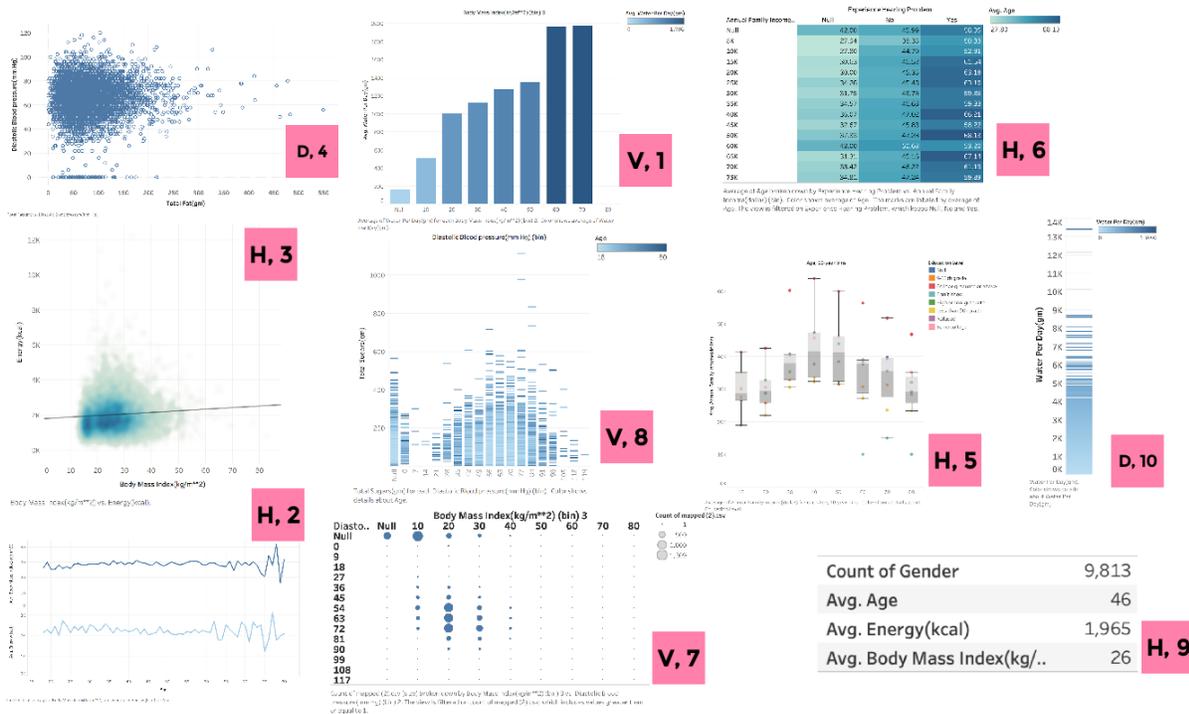


Figure 1: The gallery of visualization recommendations provided to our study participants. While we did not share the origin of these visualizations with our participants, for illustration in this paper we have labelled the chart provenance: those marked “H” were generated by pilot participants as part of their ideation sessions, those labelled “V” were recommended by Voyager [47], and those labelled “D” were recommended by Dziban [29]. We also include the overall rank of the recommendation, according to our participants. Overall rank is determined by summing reverse ranking data grouped by the visualization ID, then sorting in descending order.

visualization preferences. Though this observation could hint at the need to implement scope-reduction (in terms of attributes) within visualization recommendation systems, for logistical purposes, we decided to present just a subset that was approved by our pilot participants.

We also reduced the number of sketches we asked participants to create in the Visualization Ideation Phase of the study. We made this decision based on study length as observed through our pilot (see subsection 4.3).

4.5.2 Limitations of our Participant Pool. Though visualization recommendation systems are of broad interest for analysts across many domains, we recruited participants specifically with experience in analyzing public health data, resulting in a relatively specified participant pool. Our findings may not necessarily translate to other domains. However, given the importance of visualization in communicating public health information, and the richness of our selected dataset (NHANES), we believe our experiment still provides meaningful insights for the visualization community. We also believe that our decision to focus on participants with particular domain expertise would serve to highlight discrepancies in their design values compared to the relatively domain-agnostic priorities of algorithmic visualization recommendation systems.

We focused on participants with domain expertise who all expressed at least some familiarity with the sort of survey data that we used in our experiment. We expect that entirely different design challenges and trade-offs would result among participant pools with varying levels of familiarity with the source data.

4.5.3 Limitations of a Remote Study. Due to limitations imposed by the COVID-19 pandemic, we conducted all interviews remotely, with telecommunication software (Zoom) and an online sketching tool (Google Jamboard). Given that participants can vary widely in their experience with online tools, this mode of administering the experiment likely impacts our findings. That being said, our participants were able to successfully complete the study, and share meaningful visualization recommendations. When studies can be safely conducted in-person again, an in-person study would enable us to administer the study without the limitations imposed by our software.

4.5.4 Limitations of the Mediation Approach. As observed in prior studies [11, 41], the human mediation (or “Wizard-of-Oz”) process can sometimes lead to confusion on the part of participants, given the interplay between a participant’s intent, communication of this intent to the “wizard,” and what can be feasibly designed by the participant using real data. Many times, participants were surprised

to see the result of the visualization after their initial sketch. We observed that most of these surprises were due to a mismatch in expectations between the user's understanding of the data, and what the data actually supported.

Furthermore, we acknowledge that since the participant's intent is being interpreted by the mediator, our recreations of participants' sketches may not perfectly match their expectations. The creation of visualization recommendations under longer time frames, with deeper analyses of the data, and with additional rounds of feedback we suspect would generate categorically different sorts of visualization recommendations, but this remains an area of future study.

5 ANALYSIS

To reiterate, the main question driving our research is: *what do analysts value in visualization recommendations, and are these values reflected in the design of current visualization recommendation systems?* To answer this question, we collected study data regarding participants' design values and preferences as they created their own visualization recommendations and ranked visualizations from a pre-defined gallery of recommendations. First, we describe how we qualitatively and quantitatively analyzed this data to shed light on the perspectives and design values of analysts from the public health sector when recommending visualizations. Then, we present our analysis results, organized around the research questions listed in section 3. All of our analysis code and results are shared in our open-source repository on OSF: <https://osf.io/xeub3/>.

5.1 Qualitative Analysis Methods

We adopted a grounded theory-inspired approach for qualitatively analyzing our study data. Given our collected data (transcriptions, experimenter notes, participant sketches, and participant rankings), we used emergent coding to identify common themes in how participants design and rank visualization recommendations. In this way, we could avoid making assumptions a priori about our participants and thereby maximize our ability to observe a broad array of participant preferences. As part of our coding process, we explored potential themes regarding the visual forms of the visualizations (e.g., what encodings or visualization types our participants selected) as well as each participant's analytical intent (e.g., the types of statistical patterns our participants wanted to observe) and semantic context (e.g., the domain knowledge applied to value one visualization over another).

As a starting point, four members of the team independently went through differing subsets of at least 6 participants of interview data and created a set of themes and codes based off their assigned participants. Assignments included overlapping subsets, to not only improve saturation but also to validate that synthesized ideas were consistent for the same participant. After all sets of notes, themes, and codes were completed, one member of the team compiled an "aggregate" codebook based off the themes derived from each set. Similar themes were merged, grouping together notable quotes from participants, while broader themes were split into additional themes. For example, themes that only appeared in one set were often grouped under a broader theme. We included a brief explanation for what each code represented. We then reviewed and revised this

codebook to ensure that they consistently captured a useful range of ideas and experiences expressed by our participants.

5.2 Quantitative Analysis Methods

Our general approach to quantitatively analyzing our study data was to programmatically extract the attribute selection, data transformation, and encoding selection choices made by participants using the Vega-Lite specifications derived from each visualization observed in our interviews. Given these extracted parameters, we counted observations of each choice, such as the total participants that created visualization recommendations that use bars as the mark type or the total participants that include color encodings in their created visualizations. Here, we describe how various parameters were extracted from our study data for analysis.

5.2.1 Analyzing Recommendations Created By Participants. We translated each visualization recommendation created by participants into a corresponding visualization specification in Vega-Lite [38]. We created a series of Python scripts to analyze these specifications to extract different features, such as which attributes were selected from the NHANES dataset, what data transformations were applied (e.g., binning and aggregation), and what mark types and encoding channels were chosen to visualize these attributes. We supplemented this automatic information with manual assessments of certain properties of the designs (e.g., which variable in a bivariate chart was the independent variable and which was the dependent variable).

5.2.2 Analyzing Participants' Rankings of Existing Recommendations. Participants' rankings of visualizations from the pre-defined recommendation gallery were stored as a single relational table that recorded which 5 visualizations were selected from the gallery and each participant's ordering of their top five recommendations. Each visualization from the gallery was also translated into a corresponding Vega-Lite specification, with a note capturing recommendation source (human-generated, or auto-generated by Voyager or Dziban). We used the Vega-Lite specifications and the rankings table as inputs to scripts that extract the attributes, statistical patterns, visualization types, and encodings participants preferred from the gallery.

5.3 What do Participants Prioritize When Creating Their Own Visualizations?

In this section, we address question **RQ1** from section 3: *What characteristics of a visualization design do analysts prioritize when recommending them to colleagues?*

5.3.1 Participants created simple visualizations with a small set of intended takeaways or messages. In terms of the total attributes involved in each sketch, the vast majority of sketches include only two attributes from NHANES (37 out of 53). A small subset includes three attributes (15 out of 53), and only one example incorporated five attributes. Most recommendations appear to favor bivariate relationships between data attributes, though multiple calculations were sometimes applied to individual attributes, producing more than one encoding for this attribute within a single visualization recommendation.

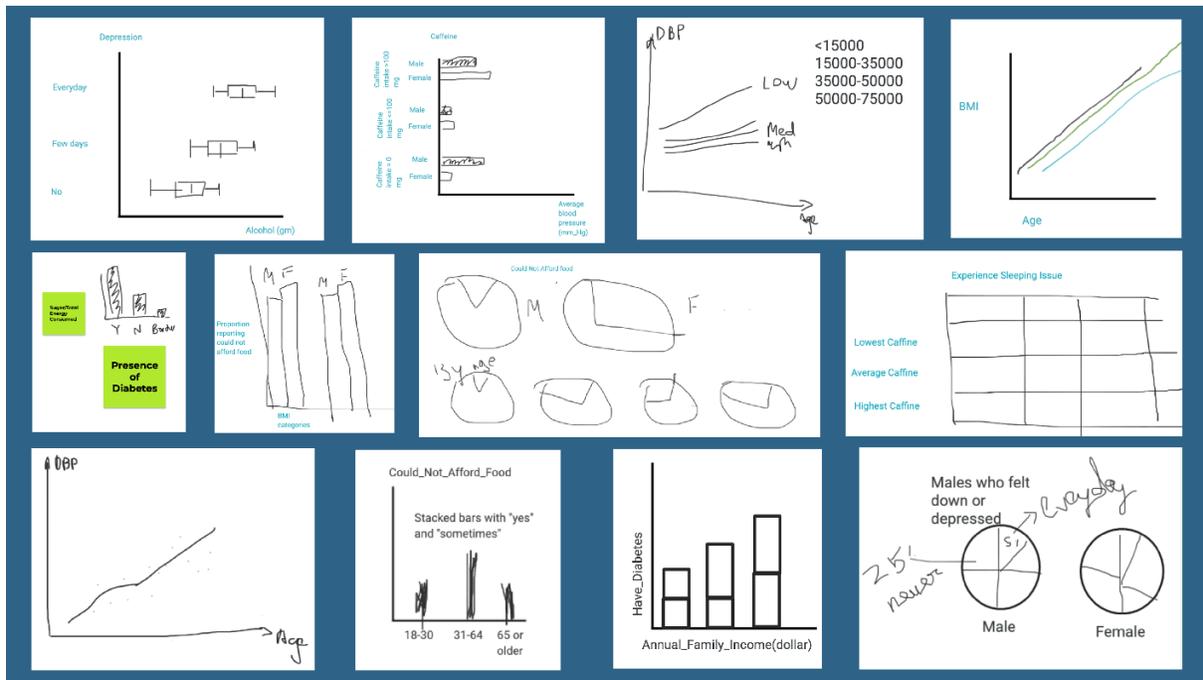


Figure 2: A gallery of recommendations sketched by participants.

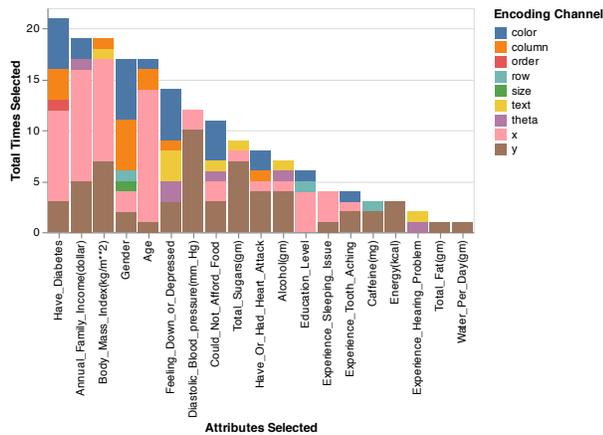


Figure 3: The frequency with which individual attributes from NHANES were used in participants' sketched visualization recommendations from the Visualization Ideation Phase, colored by how each attribute was visually encoded within each sketch. Positional encodings (such as x, y, row or column) were the modal ways that fields were mapped to visual channels.

We did not name Tableau as the system used in our interview protocol, nor did we train or inform our participants on the specific capabilities of Tableau. Our participants were not given any restrictions on the visualizations they could request. Nevertheless, participants eschewed “xenographics” [24] and stuck to common

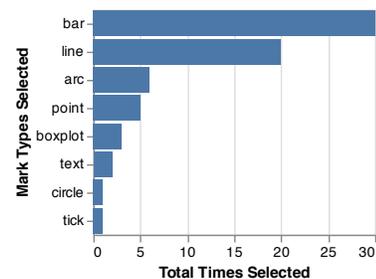


Figure 4: Visualization Ideation: Mark Representations

visual forms like bar charts, line charts, pie charts, and scatterplots (see Figure 4). We did not encounter any visualizations that could not be converted to the intentionally limited Vega-Lite [38] specification language. While participants were often aware of more complex ways of showing the data, ease of interpretation was a frequent motivator for keeping things simple. For example, P13 said “a lot of my clients aren't really data-savvy. So I use a lot of bar charts, because they are pretty easy for them to understand... if you get too crazy or too fancy, then I lose them.”

To further simplify their visualizations, participants generally applied aggregation to a dependent variable in their sketches to highlight distributions, trends, or patterns in the data. Examples include a boxplot visualization sketched by P1 to show the distribution of alcohol consumption on the x-axis (dependent variable) grouped by whether respondents felt down or depressed on the y-axis (independent variable), or a bar chart sketched by P15 showing

the percentage of respondents with diabetes on the y-axis (dependent variable) grouped by defined income classes on the x-axis (independent variable). Participants would further streamline the narrative of their visualizations by hiding or removing ambiguous records. For example, at least 9 participants wanted to de-clutter their visualizations by filtering out values such as “Don’t Know”, “Refused to answer”, and nulls. **P12** tells the visualization mediator, “Let’s filter those out and just put it in a note. Yeah, I think three [categories] in a dataset this large might be distracting to include in the graph.”

Participants also wanted clear take-aways or other action items for the intended audience, and would express disappointment when the visualization they sketched did not produce the expected clear trend. For example, when discussing all of their generated plots, **P1** was concerned that “I don’t think they’re very actionable right now.” Similarly, when discussing a potential chart of BMI, **P2** said “I probably wouldn’t send this to anybody because it looks the same across education levels, so that clearly doesn’t matter much.”

Finally, participants would use annotations or other design tweaks to highlight the intended message. **P11**, looking over a messy scatterplot, said “I... added a trendline to indicate some kind of trend so it’s easier for the audience, especially a policymaker, to see the key message.” When deciding whether or not to aggregate their data, **P8** remarked “If there’s a clear trend, then that, to me, is more meaningful than on a scatter plot. That [scatter plot] doesn’t really show anything.” **P8** also used titles to indicate intended messages (e.g. “No trend between Diabetes status and Daily Alcohol intake” and “Higher Frequency of Toothaches is Associated with Higher Daily Sugar Intakes”), a design choice also undertaken by **P13**, who remarked “I also tend to put in some kind of interpretive title or subtitle, so I tell them what they’re seeing. So they don’t have to guess.”

5.3.2 Participants focused on the visual appeal and legibility of their charts. Participants would often perform several rounds of iteration with the wizard in the ideation task in order to improve the aesthetics and legibility of the final charts. **P8** mentions of their bar chart: “I like putting outlines around the bars, so they stand out a little bit more... I also usually take off the lines going across⁴, because if you have the numbers on top of the bars, then you don’t need the lines going across like that.”) Some participants, like **P18**, conflated alternative mark types with aesthetic adjustments, mentioning: “I would keep picking at these... for a long time to keep refining them and refining them... And sometimes, I’ll make a figure as a bar chart, and then I’ll just do a stacked bar... And then I’ll try it as a pie chart. And I’ll look at them side by side and decide,” an indication that iteration on visualization designs (including tweaks such as encoding channels, mark types and even miscellaneous adjustment) serves an important part in gaining fuller confidence in their presented visualization recommendations. The lack of ability to iterate on or control the design aspects of visualization designs may have contributed towards participant’s reluctance to rank exterior recommendations more highly than their own creations (see subsection 5.5).

5.3.3 Participants focused on data attributes connected with existing contexts and shared expectations. Figure 3 lists all of the attributes that participants used in their visualization sketches from

⁴Referring to the default background gridlines produced by Tableau Desktop

the ideation section of our study (see subsection 4.4). 19 out of 20 attributes were used. The only attribute that participants did not use in their sketches was `Used_Marijuana_Or_Hashish`. For the attributes used in multiple sketches, we found that `Age` and `Education_Level` were the only attributes used exclusively as independent variables across all sketches. Similarly, `Total_Sugars` was the only attribute to be used exclusively as a dependent variable. Otherwise, attributes appeared both as independent and dependent variables within visualizations. For example, a participant might perform a breakdown of two demographic variables in either order (say, `Have_Diabetes` by `Age`, or `Age` by `Have_Diabetes`), placing the variable in a dependent or independent role depending on their specific analytic context (e.g. “Are older people more likely to have diabetes” or “Are people with diabetes more likely to be elderly”).

That being said, most demographic measures were consistently used as *independent* variables in participants’ sketches, such as `Annual_Family_Income` (used as the independent variable in 12 of 16 sketches including the attribute), `Gender` (13 of 14), `Age` (11 of 11), and `Education_Level` (4 of 4). Most health measures were generally treated as *dependent* variables, such as `Body_Mass_Index` (8 of 14 sketches including the attribute), `Have_Diabetes` (7 of 11), `Total_Sugars` (8 of 8), or `Diastolic_Blood_Pressure` (5 of 7). Most sketches involved one demographic measure and one health measure, or two health measures. The assignment of causal roles for these variables was not arbitrary. From the rationales provided by participants over the course of the ideation exercise, participants mentioned the need to identify “drivers” and “trends” (**P1**) connected to health or health outcomes. As per **P3**, “I thought about more social determinants of health, and what would be affecting it” and **P13** remarks: “I’m looking for things that relate more to getting to outcomes, whereas some of these seem more like physiological relationships or demographic relationships.”

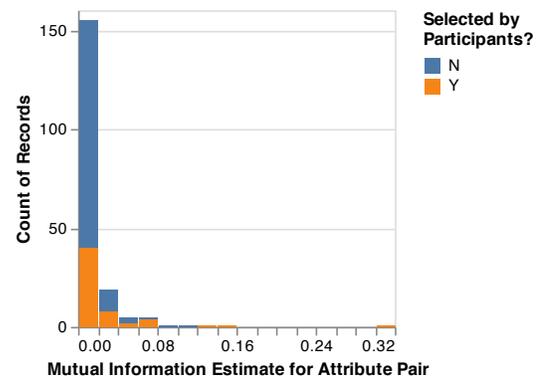


Figure 5: Counts for all possible pairs of attributes, binned by estimated mutual information value (similar to correlation) [6]. A higher mutual information estimate indicates a stronger correlation for the given pair of attributes. The pairings that were used in participants’ sketches are colored orange.

The pairing of independent and dependent variables were typically made in advance of seeing the rendered results, suggesting

that participants utilized prior knowledge and/or experiences to form hypotheses regarding potential correlations and causal relationships. That is, rather than looking for highly correlated fields and surfacing them to users, participants selected fields where the relationship was perceived as relevant or important, regardless of the resulting statistical correlation or (lack of) visible trend. For example, **P8**, upon seeing the results of a bar chart showing alcohol consumption broken down by diabetes status, remarked *“It doesn’t show me anything, should I change it? But I mean, that’s sometimes good to know to see that there isn’t really a trend between two things.”*

To confirm the disconnect between statistical correlations and the trends our participants chose to surface, we estimated a mutual information value for all possible pairs of attributes in our *NHANES* dataset (similar to a correlation score) [6]. Then, we identified which pairings were actually observed in participants’ sketches, which are colored orange in Figure 5. A higher correlation between the values in the selected attributes translates to a higher mutual information estimate. We observed a slight preference for higher ranking pairs, however, pairings span the full spectrum of observed mutual information scores.

Providing recommendations that were linked thematically or otherwise mutually supportive was also an important factor for participants, who preferred to keep a common theme among their generated visualizations. For example, **P11** (interest in Diastolic_Blood_Pressure on the y-axis for all three plots), **P12** (interest in Have_or_Had_Heart_Attack in all three plots), **P13** (interest in Have_Diabetes for all three plots), and **P3** (interest in Total_Sugars and Total_Energy in two plots) all preferred similar data attributes and/or plot types when sketching their three recommendations in order to create a synergistic, mutually supportive dashboard.

5.4 What do Participants Prioritize When Ranking Pre-Defined Visualizations?

In this section, we address question **RQ2**: *What do analysts prioritize when evaluating visualization recommendations from other sources?*

5.4.1 Participants preferred simplicity over information overload. When presented with ten visualizations, participants vastly preferred simple graphs that could be interpreted easily. While ranking the visualizations, **P1** mentions *“I think some of the graphs had a little more information than what might have been easily digestible.”* and similarly **P3** ranked the visualizations *“mostly by simplicity: which ones were quickest and easiest to follow?”* Participants also commented on the ambiguous data present in various visualizations in the gallery, also discussed in subsection 5.3.3. **P13** mentions *“Like some of the things just confuse your data a little bit. Like we still have nulls here. We have ‘Don’t Knows’, we have ‘Refused’, so like I’d probably take those out if I was doing this.”*

As with the ideation task (subsection 5.3.3), desire for simplicity was also reflected in a preference for visualizations with clear takeaways. **P4** reports *“trying to prioritize one set... where the trends jumped out a little bit more and were easier to look at to get what it’s trying to tell you at face value. That was the number one thing.”* When presented with a visualization without a “clear trend,” **P9** argues *“If it does not have a lot of variation (in the graph), then why do we do it this way?”* These sentiments echo participants’ preferences

from the ideation task, where many participants would scrap an idea if it did not lead to a visualization with positive results.

Our participants valued simplicity (and avoided complexity) in their recommendations for a variety of reasons:

- (1) The participant was not able to understand a few of the relatively complex visualizations and would not recommend them to others. For example, **P5** had difficulty interpreting one graph and said *“like this, I find really hard to understand. To me, you have to look at too many things and then make sense of what it’s telling you.”*
- (2) The participant determined that the complexity of a visualization would hinder its effectiveness for communicating the story to the target audience (policymakers). **P15** explains *“And, some of these I think are just so complicated for people that understand what they are so I’m not choosing them. I’m also just trying to think about what would be easy for my client to show and what people would get.”* Referring to the same visualization (Visualization ID A), **P6** liked it because *“...it may be a little less fancy or something like that, but it communicates the information pretty easily, which is always good.”*

Our quantitative analyses reinforced our qualitative findings here. In terms of encoding count, we found that participants preferred recommendations representing bi-variate relationships over more complex visualizations: 48.9% of participants’ selected gallery visualizations contained exactly two attributes, 41.1% of participants’ selected gallery visualizations contained three attributes, sharply dropping off to 6.7% for gallery visualizations containing four attributes.

5.4.2 Recommendations were often considered inspiration for later exploration but not necessarily a final product. Through ranking the pre-defined gallery of recommendations, participants discussed the things they liked and disliked for each visualization. Several considered the visualizations interesting which helped spark ideas that they had not come up with while generating their own visualizations. **P5** was inspired by one area chart, and wanted to explore it further. *“I think it could be super cool. But as an area chart, I guess I’d have to play around with the data to see.”* **P1** noted that they’d like to dive into the recommended visualizations to look at subsets of the graph: *“I think I would like take cuts, like look at a subset of people who are feeling down and depressed, and then really analyze that.”*

These findings suggest that the visualization recommendation process should be treated as a dynamic, evolving dialogue between the user and the system, where recommendations can spur new ideas but not necessarily substitute for the user’s own visualization design work.

5.5 How Do Participant Recommendations Compare With Those From Other Sources?

In this section, we address question **RQ3**: *How do the recommendations made by analysts align with those created from other sources in terms of visual form or analytical purpose?*

5.5.1 Participants preferred their own creations to recommendations from others. Only 6 out of 18 participants picked gallery visualizations to replace their own created visualizations. Our observations suggest several potential reasons that together point to a participants' perceived difficulty parsing the visualization, including: aesthetics & visual formatting (“*This one is sort of hard to read since the x-axis isn't labeled and that just sort of hard figuring out what it's saying.*” -P12), and simplicity or clarity (“*I rarely make a visual for a technical audience, and so I look at these charts and I'm like, my client wouldn't understand this.*” -P5). These factors were occasionally valued over even analytic utility or accuracy, e.g., “*Some of these graphs, like A or D, might actually be describing the dataset better than the ones I selected. But because it's so much harder to read... that's not something I would show a client. That's why I didn't even go further than an initial look.*” (P19). Together, these findings suggest that analysts are consistent in their recommendation preferences—in this case, prioritizing clear and simple narratives—regardless of whether they are creating new recommendations or ranking existing ones.

5.5.2 Participants exhibited no strong preferences for recommendations generated by other humans versus those generated by recommender systems. In terms of recommendation sources (e.g., human versus algorithm), we found that when selecting from the gallery, participants had a slight preference for human-curated recommendations (48 selections) over those generated using Voyager and Dziban (24 and 18 selections, respectively – 42 total). However, we note that our piloting process had, in a way, already pre-filtered algorithmically generated recommendations (see subsection 4.3) based on interest. We are therefore hesitant to use our quantitative results to point to a strong pattern of preference for a particular source of recommendation, beyond pointing out that there was no extreme bias towards one source of recommendations or another.

5.6 Reflections from the Interviewers

Our original preregistration focused on determining how our participants perceived the recommendation design process, but did not consider the perspective of the interviewers (the mediator and the wizard). In this section, the interviewers in the experiment reflect on their experiences and observations from the interview transcripts. We deviate from our original preregistration to include these perspectives as we believe they provide richer nuance surrounding opportunities and challenges within the visualization recommendation process.

We highlight notable misalignments in expectations that we observed throughout our study, clustered according to three recurring themes: uncertainties about the data from the user's perspective, uncertainties about user intent from the wizard's perspective, and the inability to generate the visualization according to the original user intent. These misalignments led to interesting compromises in the end-visualization, as detailed below.

5.6.1 Assumptions and Uncertainties about Data. On several occasions, the wizard interfered with the participant's design process by reminding them of characteristics of the data. However, on one occasion, this interruption spun into the participant selecting a radically different choice in data attributes.

Wizard: For the second line plot where you had income by BMI... One thing to keep in mind is that income isn't actually a continuous variable; it's discrete with a size of 5k.

P6: Thanks for pointing that out, I did not notice that. Let's do something else then... Let's just ask a curious question. Instead of the BMI and salary, let's do caffeine, by age and by gender.

While in some instances (such as the above) this new awareness or clarification about the data caused the participant to change plans entirely, in the vast majority of scenarios, the wizard and the participant worked to derive meaningful compromise on the data attributes once the data was clarified. For example, this snippet from the P18 transcript shows an incremental back-and-forth which fleshed out a visualization design, as gaps and uncertainties in knowledge of the dataset were resolved between the wizard and the participant:

Wizard: Just a quick note that for the first graph you sketched, “Could Not Afford Food” actually can contain three different values.

P18: Thank you for pointing that out... I shouldn't have assumed it was [boolean].. Can you do it as like a stacked bar [instead of a standard bar chart]?

In addition to highlighting the value of mutual communication and iteration when creating visualizations, these interactions also suggest the utility in presenting users with useful summaries or alerts about the data available to them in order to orient them to their data.

Given the shape and size of our chosen dataset, generating many standard visualizations required an explicit choice of aggregation. Our own participants often did not explicitly qualify an aggregation method or took it for granted that the system would use their intended default. This likely extends to when they interact with drag-and-drop interfaces such as Tableau. In the experiment, the wizard had to thus frequently prompt for explicit aggregation instructions, such as in the P16 transcript:

Wizard [referring to bars, after P16 finished sketching]: So would it be like... the average sugar intake?

P16: Mmm.. Yeah! You could do that [averaging] for any of the health outcomes and see your relationship with annual family income.

...

Wizard [referring to a second sketch]: For the y-axis, did you want to zero in on one attribute in BMI-sugar-caffeine? Should we just pick one?

P16. Sure. BMI, although that seems like probably well-established already in a policy.

Wizard: Okay, so for now, I'll do BMI. Do you have an aggregation measure you'd like to pick? What would you like to do? I'd like to remind you that in the first graph, we decided to average the total sugar intake.

P16: Again, to average the amount.

This uncertainty around aggregation type has been observed not just in creators but also readers of visualizations [18] and suggests

that choosing the wrong “default” could produce confusing or even misleading recommendations.

Lastly, the participants often asked questions that required deep knowledge of either the domain or the source data. As the wizard and the mediator were not public health experts themselves, these questions occasionally led to acknowledgments of uncertainty:

P11: For the diabetes question, how was that question worded?

Mediator: Regarding the exact question, I can't recall off the top of my mind, but the user can answer either "Yes", "No", and "Borderline".

P11: Okay, so the question didn't mention anything about diagnosis?

Mediator: I'm not 100% sure.

Some of these details about the dataset matter and affect how a participant would generate a visualization. We did not always have the answers to some of these questions. In a similar way, many visualization systems do not have all of the semantic details about a dataset it might visualize for a user, though this would streamline the design process by improving the user's literacy on the dataset and improving their ability to find appropriate visualization forms for their data.

5.6.2 Uncertainties About User Intent. Analogously to the participant's uncertainties about the data, there were sometimes uncertainties on the wizard's perspective about user intent. This sometimes led to scenarios that even after a couple clarifying questions, we were unable to 100% follow what a user expected. This was often compounded when the participant would provide vague visuals.

P5: So I don't know how doable this is, I would love to, and there may not be time for this for like what you have planned for today... So I'm thinking like, if we had more time, it would be cool to group education level differently. Right. So like, anyway, so like, less than ninth grade, and ninth to 11th grade could be grouped as one, but obviously, we don't need to do all that grouping now. But like that grouping would just have two dots on its line. Right, it would show the x axis which would be the spectrum of no to every day, like no over at the zero, what up to essentially, the education grouping would have two dots on its line, as opposed to right now it has 10 you see what I'm saying? Yeah, it's far more simple. Or we can move on to the other one.

Wizard: Okay, let's keep this for now. And then if we have time, we can go back and do these adjustments.

There was also at least one case in which the wizard was unable to meet the participant's needs with an alternative visualization design.

Wizard: Oh, yeah, that um, that is up to you if you have a different way that you want to visualize it. Since the table is a little bit tricky.

P20: So I think then we can go for line graph also.

Wizard: Okay. Yeah, so I think for for this one, we have a line graph, it can't be done since there are categories for the body, since it's because of the different categories. So let me show you. Like we can have this bar chart with the annual family income versus each of the different types of different types of BMI. Yeah, so is this an alternate solution?

P20: Okay. This is still not quite what I am looking for.

Interviewer: Right, I guess, if that's the case, we can still mark that you still want the table? Because I realized we're like overtime a little bit... Yeah, we can definitely take notes on that and maybe mark that down in our notes.

These moments ultimately affected the wizard's ability to address or tweak user intent in the design. While having a *wizard* in the loop was in some sense detrimental for full expressiveness (since the wizard may have had imperfect knowledge of the system), it also afforded us fuller *transparency* in communicating both system and wizard limitations to users, as well as *mindfulness* on the part of the participants of potential designs that required significant time or effort to “get right.” That being said, it is presently difficult to provide the same expressiveness a user might have with pen-and-paper to a visualization system, especially when factors such as data, domain, and technical know-how exist. Compromises must occur, which is why we detail these moments in our interviews here.

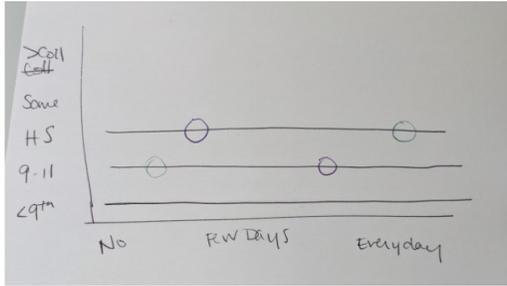
5.6.3 Constraints: System and Wizard Ability. While the wizard tried to maintain the original user intent as intact as possible, this wasn't always possible due to constraints on the system (Tableau) or the wizard's ability to use the system to create a visualization isomorphic to the participant's initial design, in the time allotted. Therefore, while we always recorded the original sketch and design intention of the participant, on several occasions, the wizard and the user were forced to produce a compromised design, or abandon potential design ideas, sometimes due to the request for more advanced analytical queries.

Wizard: So let me try to understand what you were looking for. So you're looking for a proportion graph per income level, where the number essentially represented is the numerator [family income] divided by the total number of respondents across all family incomes?

P9: No. So like, let's say there are 100 respondents in the 5k category. 100 males and 100 females. I want to see of those 100 males, what proportion reported feeling down or depressed? Maybe it's only 10%? And then among females, what proportion reported feeling down or depressed?

Though we tried to limit drastic interference (e.g. “here's a radically different visualization design idea you could consider”), sometimes compromises due to technical skill with the visualization system had to be made.

P9: I'm just trying to imagine what that could look like, besides doing like side-by-side bars for the



(a) P5 sketch. P5 notes that nodes are colored according to the dataset’s gender attribute.



(b) Wizard’s interpretation of P5 sketch

Figure 6: P5’s initial design idea (Figure 6a) and the wizard’s interpretation of the design (Figure 6b). These visualizations have slight differences, including filtered values on both axes, as well as more structural differences such the data grid inherent in Figure 6b. For this sketch, the wizard was unclear on the design intent of P5, specifically regarding the mark types, which ultimately led to their inability to specify this in Tableau.

males and the females. Do you have any suggestions?

Wizard: Well, we could just separate the axes and make one axis female and one male.

P9: Okay, yeah I like that. So female versus male.

The examples in this section highlight that the wizard was not a real-time translator of a participant’s visualization idea into a final design, but had to navigate their own technical skill with Tableau, along with the uncertainties and ambiguities regarding user intent and data. For all sketches, the wizard was able to generate a design compromise that the participant accepted for the sake of “moving on.” However, the scenario demonstrated in the snippet with P9 may not translate to a real-world scenario where certain analytic functions are necessarily to be expressed in the visualization. Sometimes, the request of advanced technical functions on data attributes led to more creative, bespoke visualizations in which the wizard, though decently familiar with Tableau functions, was not technically prepared to reconstruct during the interview, hampering the participant’s visualizations. Although we were mostly able to come to a compromise with the participant’s visualization design, it becomes more difficult to manage their expectations for visualizations involving complex transforms, which may not be easily replaceable with a simpler one, when they are requested.

6 DISCUSSION

While there were many themes and potential vignettes related to our data, we focus our discussion on three recurring themes encountered in our analysis. These themes represent *design values* or *priorities* that recurred across participants and across our experimental tasks: **simplicity**, **relevance**, and **interest**. We use this section to explicate these design values, with a focus on how these

values might inform current or future designers of visualization recommendation systems. We reiterate that our choice of participants and methods results in knowledge that is positional, and exists in the context of our pool of public health researchers and analysts. What counts as “simple, relevant, and interesting” for them may not necessarily hold true for other domains or target audiences.

6.1 Visualizations Should Be Simple

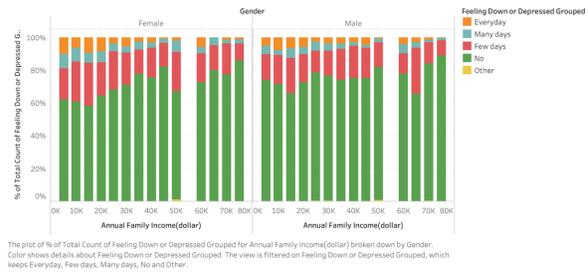
As discussed in section 5, participants repeatedly preferred *simple* visualizations with minimal attributes both when creating and ranking visualizations. This notion of simplicity extends to both the visual *design* of the visualizations as well as the *messages* of these visualizations. In terms of *design*, this simplicity is associated with a preference for familiar visualization designs (bar charts, line charts, and scatterplots) encoding only one or two variables at a time. In terms of *messages*, this simplicity was reflected in a preference for clear or otherwise unambiguous statistical patterns (such as expected trends or relationships between two variables), with filtering and aggregating used to remove extraneous information from a chart.

Our participants provided several rationales for preferring simplicity, but two of the most relevant for designers of recommendation systems were the desire to present clear takeaways for an audience with variable experience interpreting complex or unfamiliar charts, and the sense that visualization recommendations should serve as a starting point for future analyses, rather than an ending point. From these design values, we propose two guidelines for designers of recommendation systems. Namely, recommendation systems should:

- (1) **produce charts from common genres with bounded visual complexity.** A visualization recommendation system, especially one for the purpose of orienting or familiarizing



(a) P9 sketch.



(b) The wizard's interpretation of P9 sketch

Figure 7: Participant P9's initial design idea (Figure 7a) and the wizard's interpretation of the design (Figure 7b). Despite the apparent simplicity of Figure 7a, P9 paired it with more advanced analytic functions that were too technically challenging for the wizard to implement in the time allotted. The design compromise is shown in Figure 7b.

a user with a dataset, is not the place to try esoteric or unfamiliar charts. Similarly, rather than attempting to include every potential relevant dimension in one design, the recommendation system should value and present potentially greater sets of univariate or bivariate charts.

- (2) **offer opportunities for refinement and extension.** Especially in the context of our suggestion above, a single “simple” visualization may not be sufficient to meet the analytical goals of all users. Rather than increase the complexity of the recommended chart to cover all potential goals, systems should allow users to modify or extend recommendations within more traditional exploratory data analysis workflows. Rather than an “automatic insight,” a recommendation might be better thought of as an interesting question or hypothesis to verify in more detail or with more nuance after the fact.

Note that a person's perception of simplicity can also be influenced by their prior experiences with certain encoding channels, which is dictated in part by their discipline. For example, domain experts may prefer to use established encoding conventions [2] from within their discipline, even though these encodings may technically be more complex or less legible [8]. Additional studies could prove beneficial in distinguishing domain-specific guidelines from design principles that could apply to any type of visualization recommendation system.

6.2 Visualizations Should Be Relevant

As mentioned in section 2, many visualization recommendation systems are built to be agnostic to the data domain, providing recommendations purely based on design guidelines, statistical features, or the syntactic structure of the data. Our participants, however, were more motivated by the *anticipated needs* of their audience and their own *structural and semantic* understanding of the data. As with simplicity above, we acknowledge that *relevance* is similarly polysemic and nuanced as a concept.

For instance, participants often zeroed in on specific variables *a priori*, excluding variables they felt to be irrelevant to their likely

audience or unrelated to another pre-selected variable of interest. These choices represent implicit or explicit *hypotheses* about the relations in the data, as well as implicit or explicit *causal modeling* of how one factor might influence another. Providing evidence for or against these implicit hypotheses resulted in a common recipe for bivariate visualizations: a demographic factor cast as an independent variable, with a health outcome measure cast as a dependent variable.

The centrality of domain relevance leads us to make the following recommendations for designers of visualization recommendation systems. Namely, systems should:

- (1) **incorporate domain contexts and field relationships.** This could be as simple as giving common fields like *Age* or *Gender* fixed roles (say, as independent variables) in recommended charts, or as complex as attempting to infer or solicit causal graphs, functional dependencies, or other field relationships in particular datasets. This recommendation also suggests that the promise of a *universal* “auto-insight” system is potentially misguided, or at least a lofty dream, and systems seeking to present relevant data may need to be tailored to specific domains or users.
- (2) **allow users to specify intent or their analytical goals.** The exact same visualization may be alternatively useful or irrelevant across different users depending on their interests or prior assumptions. Allowing users to specify particular fields of interest or even their hypotheses and assumptions (see below) would allow systems to be more adaptable and retain relevancy across users or analytics sessions. While Q&A systems (see subsection 2.1.3) have begun to incorporate models of intent into their design, we believe that this sort of contextual or ancillary semantic information is critical across all sorts of recommendation systems.

6.3 Visualizations Should Be Interesting

The last design value we examine in the context of our findings is that of preferring *interesting* visualizations over visualizations that

gave the viewer (or the audience) no new and useful information. This need for a visualization to show *something* (a trend, a skew, or just the more general notion of a “takeaway”) has potentially troubling implications for system designers.

A common pattern was for participants to abandon (in our ideation task) or downweight (in our selection task) visualizations without definitive patterns and trends. This includes visualizations showing negative or null results with respect to an assumed trend or difference between or among groups. P8’s verbal quandary (see subsection 5.3.1) over whether to change a visualization that “doesn’t show me anything” was an exception in this respect. In other words, critiques of auto-insight systems (see subsection 2.1.1) based on their potential to highlight dramatic but spurious effects may be similarly applied to human recommenders. While we therefore *could* suggest that recommendation systems present only the most salient of positive results, as an alternative we present the following suggestions. Namely, that systems should:

- (1) **allow users to specify explicit hypotheses or assumptions.** A chart showing absolutely no correlations or clear patterns might still be of interest to a viewer with a strong expectation or assumption about a relationship in the data. To afford these sorts of interesting findings, recommendation systems (especially auto-insight systems) should embrace what Hullman & Gelman [16] refer to as “theories of graphical inference” such that visualizations can provide direct evidence for or against inferences about the data. Kim et al. [20] in particular suggest externalizing expectations and predictions prior to looking at the data as an alternative to traditional data presentation.
- (2) **provide additional guidance to users on the robustness or importance of the visual pattern in a particular recommendation.** To prevent users from being misled by potential spurious patterns in sets of visualizations (such as those encountered and reported by the participants observed by Zraggen et al. [50]), recommendation systems may need to provide additional information on the reliability of findings (either through existing or novel metrics of reliability, or through alerts or warnings of potential threats to validity [33]).

7 CONCLUSION

This work explored the preferences and priorities of in-domain analysts when creating and evaluating visualization recommendations for their target audience. We base our findings off semi-structured interviews with 18 analysts in the public health sector, observing behaviors, attitudes, and perceptions they had for various components throughout the visualizations in the experiment. Our findings highlight that these users overwhelmingly value **simplicity, relevancy, and analytic interest** both when creating their own visualization designs and when evaluating other visualization recommendations. Participants either demonstrated an understanding of a visualization and the potential story or insight stemming from it, or were able to formulate targeted questions to any uncertainties with the attributes in the visualization. Furthermore, certain data attributes (demographic and health outcomes) were frequently paired in visualization designs, indicating natural priors and biases in semantic

knowledge in the dataset. Lastly, we find that participants more likely engage with visualizations showing seemingly “positive” results with perceptively clearer patterns or trends. Based on our findings, we suggest various design possibilities for visualization recommendation designers that could better aid in data exploration workflows.

ACKNOWLEDGMENTS

We thank all the reviewers, study participants, and members of both the Human-Computer Interaction Lab and the Battle Data Lab for their valuable feedback.

REFERENCES

- [1] Inc. Amazon Web Services. 2021. Amazon QuickSight - Business Intelligence Service - Amazon Web Services. <https://aws.amazon.com/quicksight/>
- [2] Matthew Brehmer, Jocelyn Ng, Kevin Tate, and Tamara Munzner. 2016. Matches, Mismatches, and Methods: Multiple-View Workflows for Energy Portfolio Analysis. *IEEE Transactions on Visualization and Computer Graphics* 22, 1 (2016), 449–458. <https://doi.org/10.1109/TVCG.2015.2466971>
- [3] Centers for Disease Control and Prevention (CDC). National Center for Health Statistics (NCHS). 2013. National Health and Nutrition Examination Survey Data. <https://www.cdc.gov/nchs/nhanes/index.htm>.
- [4] Remco Chang, Caroline Ziemkiewicz, Tera Marie Green, and William Ribarsky. 2009. Defining Insight for Visual Analytics. *IEEE Computer Graphics and Applications* 29, 2 (2009), 14–17. <https://doi.org/10.1109/MCG.2009.22>
- [5] Michael Correll. 2019. *Ethical Dimensions of Visualization Research*. Association for Computing Machinery, New York, NY, USA, 1–13. <https://doi.org/10.1145/3290605.3300418>
- [6] Thomas M Cover and Joy A Thomas. 2012. *Elements of Information Theory*. John Wiley & Sons. <https://doi.org/10.1002/047174882X>
- [7] A. Crisan, S. Fisher, J. L. Gardy, and T. Munzner. 5555. GEViTRec: Data Reconnaissance Through Recommendation Using a Domain-Specific Visualization Prevalence Design Space. *IEEE Transactions on Visualization & Computer Graphics* 01 (aug 5555), 1–1. <https://doi.org/10.1109/TVCG.2021.3107749>
- [8] Arita Dasgupta, Jorge Poco, Bernice Rogowitz, Kyungsik Han, Enrico Bertini, and Claudio T. Silva. 2020. The Effect of Color Scales on Climate Scientists’ Objective and Subjective Performance in Spatial Data Analysis Tasks. *IEEE Transactions on Visualization and Computer Graphics* 26, 3 (2020), 1577–1591. <https://doi.org/10.1109/TVCG.2018.2876539>
- [9] Çağatay Demiralp, Peter J. Haas, Srinivasan Parthasarathy, and Tejaswini Pedapati. 2017. Foresight: Recommending Visual Insights. *Proc. VLDB Endow.* 10, 12 (aug 2017), 1937–1940. <https://doi.org/10.14778/3137765.3137813>
- [10] Victor Dibia and Çağatay Demiralp. 2019. Data2Vis: Automatic Generation of Data Visualizations Using Sequence-to-Sequence Recurrent Neural Networks. *IEEE Computer Graphics and Applications* 39, 5 (2019), 33–46. <https://doi.org/10.1109/MCG.2019.2924636>
- [11] Lars Grammel, Melanie Tory, and Margaret-Anne Storey. 2010. How Information Visualization Novices Construct Visualizations. *IEEE Transactions on Visualization and Computer Graphics* 16, 6 (2010), 943–952. <https://doi.org/10.1109/TVCG.2010.164>
- [12] Lars Grammel, Melanie Tory, and Margaret-Anne Storey. 2010. How Information Visualization Novices Construct Visualizations. *IEEE Transactions on Visualization and Computer Graphics* 16, 6 (2010), 943–952. <https://doi.org/10.1109/TVCG.2010.164>
- [13] Marti Hearst, Melanie Tory, and Vidya Setlur. 2019. Toward Interface Defaults for Vague Modifiers in Natural Language Interfaces for Visual Analysis. In *2019 IEEE Visualization Conference (VIS)*. 21–25. <https://doi.org/10.1109/VISUAL.2019.8933569>
- [14] Jeffrey Heer. 2019. Agency plus automation: Designing artificial intelligence into interactive systems. *Proceedings of the National Academy of Sciences* 116, 6 (2019), 1844–1850. <https://doi.org/10.1073/pnas.1807184115> arXiv:<https://www.pnas.org/content/116/6/1844.full.pdf>
- [15] Kevin Hu, Michiel A. Bakker, Stephen Li, Tim Kraska, and César Hidalgo. 2019. *VizML: A Machine Learning Approach to Visualization Recommendation*. Association for Computing Machinery, New York, NY, USA, 1–12. <https://doi.org/10.1145/3290605.3300358>
- [16] Jessica Hullman and Andrew Gelman. 2021. Designing for Interactive Exploratory Data Analysis Requires Theories of Graphical Inference. *Harvard Data Science Review* (30 7 2021). <https://doi.org/10.1162/99608f92.3ab8a587> <https://hdsr.mitpress.mit.edu/pub/w075glo6>.
- [17] Pawandeep Kaur and Michael Owonibi. 2017. A Review on Visualization Recommendation Strategies, Vol. 4. SCITEPRESS, 266–273. <https://doi.org/10.5220/0006175002660273>

- [18] Younghoon Kim, Michael Correll, and Jeffrey Heer. 2019. Designing Animated Transitions to Convey Aggregate Operations. *Computer Graphics Forum (Proc. EuroVis)* (2019). <http://idl.cs.washington.edu/papers/animated-aggregate-operations>
- [19] Younghoon Kim, Kanit Wongsuphasawat, Jessica Hullman, and Jeffrey Heer. 2017. GraphScape: A Model for Automated Reasoning about Visualization Similarity and Sequencing. In *Proceedings of the 2017 CHI Conference on Human Factors in Computing Systems* (Denver, Colorado, USA) (CHI '17). Association for Computing Machinery, New York, NY, USA, 2628–2638. <https://doi.org/10.1145/3025453.3025866>
- [20] Yea-Seul Kim, Katharina Reinecke, and Jessica Hullman. 2017. Explaining the Gap: Visualizing One's Predictions Improves Recall and Comprehension of Data. In *Proceedings of the 2017 CHI Conference on Human Factors in Computing Systems* (Denver, Colorado, USA) (CHI '17). Association for Computing Machinery, New York, NY, USA, 1375–1386. <https://doi.org/10.1145/3025453.3025592>
- [21] David Kirsh. 2010. Thinking with external representations. *AI & SOCIETY* 25, 4 (Nov. 2010), 441–454. <https://doi.org/10.1007/s00146-010-0272-8>
- [22] Bum chul Kwon, Brian Fisher, and Ji Soo Yi. 2011. Visual analytic roadblocks for novice investigators. In *2011 IEEE Conference on Visual Analytics Science and Technology (VAST)*. 3–11. <https://doi.org/10.1109/VAST.2011.6102435>
- [23] Heidi Lam, Enrico Bertini, Petra Isenberg, Catherine Plaisant, and Sheelagh Carpendale. 2012. Empirical Studies in Information Visualization: Seven Scenarios. *IEEE Transactions on Visualization and Computer Graphics* 18, 9 (2012), 1520–1536. <https://doi.org/10.1109/TVCG.2011.279>
- [24] Maarten Lambrechts. 2021. <https://xeno.graphics/about/>
- [25] Po-Ming Law, Alex Endert, and John Stasko. 2020. Characterizing Automated Data Insights. In *2020 IEEE Visualization Conference (VIS)*. 171–175. <https://doi.org/10.1109/VIS47514.2020.00041>
- [26] Bongshin Lee, Rubaiat Habib Kazi, and Greg Smith. 2013. SketchStory: Telling More Engaging Stories with Data through Freeform Sketching. *IEEE Transactions on Visualization and Computer Graphics* 19, 12 (2013), 2416–2425. <https://doi.org/10.1109/TVCG.2013.191>
- [27] Sukwon Lee, Sung-Hee Kim, Ya-Hsin Hung, Heidi Lam, Youn-Ah Kang, and Ji Soo Yi. 2016. How do People Make Sense of Unfamiliar Visualizations?: A Grounded Model of Novice's Information Visualization Sensemaking. *IEEE Transactions on Visualization and Computer Graphics* 22, 1 (2016), 499–508. <https://doi.org/10.1109/TVCG.2015.2467195>
- [28] Quan Li, Huanbin Lin, Chunfeng Tang, Xiguang Wei, Zhenhui Peng, Xiaojuan Ma, and Tianjian Chen. 2021. Exploring the “Double-Edged Sword” Effect of Auto-Insight Recommendation in Exploratory Data Analysis. In *CEUR Workshop Proceedings*, Vol. 2903.
- [29] Halden Lin, Dominik Moritz, and Jeffrey Heer. 2020. Dziban: Balancing Agency & Automation in Visualization Design via Anchored Recommendations. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems* (Honolulu, HI, USA) (CHI '20). Association for Computing Machinery, New York, NY, USA, 1–12. <https://doi.org/10.1145/3313831.3376880>
- [30] Jock Mackinlay. 1986. Automating the Design of Graphical Presentations of Relational Information. *ACM Trans. Graph.* 5, 2 (apr 1986), 110–141. <https://doi.org/10.1145/22949.22950>
- [31] Jock Mackinlay, Pat Hanrahan, and Chris Stolte. 2007. Show Me: Automatic Presentation for Visual Analysis. *IEEE Transactions on Visualization and Computer Graphics* 13, 6 (2007), 1137–1144. <https://doi.org/10.1109/TVCG.2007.70594>
- [32] Andrew McNutt, Anamaria Crisan, and Michael Correll. 2020. Divining Insights: Visual Analytics Through Cartomancy. In *Extended Abstracts of the 2020 CHI Conference on Human Factors in Computing Systems* (Honolulu, HI, USA) (CHI EA '20). Association for Computing Machinery, New York, NY, USA, 1–16. <https://doi.org/10.1145/3334480.3381814>
- [33] Andrew McNutt, Gordon Kindlmann, and Michael Correll. 2020. Surfacing Visualization Mirages. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems* (Honolulu, HI, USA) (CHI '20). Association for Computing Machinery, New York, NY, USA, 1–16. <https://doi.org/10.1145/3313831.3376420>
- [34] mihart. [n. d.]. Types of Insights supported by Power BI - Power BI. <https://docs.microsoft.com/en-us/power-bi/consumer/end-user-insight-types>
- [35] Dominik Moritz, Chenglong Wang, Greg L. Nelson, Halden Lin, Adam M. Smith, Bill Howe, and Jeffrey Heer. 2019. Formalizing Visualization Design Knowledge as Constraints: Actionable and Extensible Models in Draco. *IEEE Transactions on Visualization and Computer Graphics* 25, 1 (2019), 438–448. <https://doi.org/10.1109/TVCG.2018.2865240>
- [36] Evan M. Peck, Sofia E. Ayuso, and Omar El-Etr. 2019. *Data is Personal: Attitudes and Perceptions of Data Visualization in Rural Pennsylvania*. Association for Computing Machinery, New York, NY, USA, 1–12. <https://doi.org/10.1145/3290605.3300474>
- [37] Xiaoying Pu and Matthew Kay. 2018. The Garden of Forking Paths in Visualization: A Design Space for Reliable Exploratory Visual Analytics : Position Paper. In *2018 IEEE Evaluation and Beyond - Methodological Approaches for Visualization (BELIV)*. 37–45. <https://doi.org/10.1109/BELIV.2018.8634103>
- [38] Arvind Satyanarayan, Dominik Moritz, Kanit Wongsuphasawat, and Jeffrey Heer. 2017. Vega-Lite: A Grammar of Interactive Graphics. *IEEE Transactions on Visualization and Computer Graphics* 23, 1 (2017), 341–350. <https://doi.org/10.1109/TVCG.2016.2599030>
- [39] Michael Sedlmair, Miriah Meyer, and Tamara Munzner. 2012. Design Study Methodology: Reflections from the Trenches and the Stacks. *IEEE Transactions on Visualization and Computer Graphics* 18, 12 (2012), 2431–2440. <https://doi.org/10.1109/TVCG.2012.213>
- [40] Arjun Srinivasan, Steven M. Drucker, Alex Endert, and John Stasko. 2019. Augmenting Visualizations with Interactive Data Facts to Facilitate Interpretation and Communication. *IEEE Transactions on Visualization and Computer Graphics* 25, 1 (2019), 672–681. <https://doi.org/10.1109/TVCG.2018.2865145>
- [41] Melanie Tory and Vidya Setlur. 2019. Do What I Mean, Not What I Say! Design Considerations for Supporting Intent and Context in Analytical Conversation. In *2019 IEEE Conference on Visual Analytics Science and Technology (VAST)*. 93–103. <https://doi.org/10.1109/VAST47406.2019.8986918>
- [42] Barbara Tversky. 2008. Making Thought Visible. In *Proceedings of the International Workshop on Studying Design Creativity*. Springer The Netherlands.
- [43] Manasi Vartak, Sajjadur Rahman, Samuel Madden, Aditya Parameswaran, and Neoklis Polyzotis. 2015. S-span class="smallcaps SmallerCapital">eeDB: Efficient Data-Driven Visualization Recommendations to Support Visual Analytics. *Proc. VLDB Endow.* 8, 13 (sep 2015), 2182–2193. <https://doi.org/10.14778/2831360.2831371>
- [44] Jagoda Walny, Sheelagh Carpendale, Nathalie Henry Riche, Gina Venolia, and Philip Fawcett. 2011. Visual Thinking In Action: Visualizations As Used On Whiteboards. *IEEE Transactions on Visualization and Computer Graphics* 17, 12 (2011), 2508–2517. <https://doi.org/10.1109/TVCG.2011.251>
- [45] Jagoda Walny, Samuel Huron, and Sheelagh Carpendale. 2015. An exploratory study of data sketching for visual representation. *Computer Graphics Forum* 34, 3 (2015), 231–240. <https://doi.org/10.1111/cgf.12635>
- [46] Richard Wesley, Matthew Eldridge, and Pawel T. Terlecki. 2011. An Analytic Data Engine for Visualization in Tableau. In *Proceedings of the 2011 ACM SIGMOD International Conference on Management of Data* (Athens, Greece) (SIGMOD '11). Association for Computing Machinery, New York, NY, USA, 1185–1194. <https://doi.org/10.1145/1989323.1989449>
- [47] Kanit Wongsuphasawat, Zening Qu, Dominik Moritz, Riley Chang, Felix Ouk, Anushka Anand, Jock Mackinlay, Bill Howe, and Jeffrey Heer. 2017. Voyager 2: Augmenting Visual Analysis with Partial View Specifications. In *Proceedings of the 2017 CHI Conference on Human Factors in Computing Systems* (Denver, Colorado, USA) (CHI '17). Association for Computing Machinery, New York, NY, USA, 2648–2659. <https://doi.org/10.1145/3025453.3025768>
- [48] Rachael Zehrung, Astha Singhal, Michael Correll, and Leilani Battle. 2021. *Vis Ex Machina: An Analysis of Trust in Human versus Algorithmically Generated Visualization Recommendations*. Association for Computing Machinery, New York, NY, USA. <https://doi.org/10.1145/3411764.3445195>
- [49] Zehua Zeng, Phoebe Moh, Fan Du, Jane Hoffswell, Tak Yeon Lee, Sana Malik, Eunye Koh, and Leilani Battle. 2021. An Evaluation-Focused Framework for Visualization. *IEEE Transactions on Visualization and Computer Graphics* (2021), 11. to appear.
- [50] Emanuel Zgraggen, Zheguang Zhao, Robert Zeleznik, and Tim Kraska. 2018. *Investigating the Effect of the Multiple Comparisons Problem in Visual Analysis*. Association for Computing Machinery, New York, NY, USA, 1–12. <https://doi.org/10.1145/3173574.3174053>
- [51] Mengyu Zhou, Qingtao Li, Xinyi He, Yuejiang Li, Yibo Liu, Wei Ji, Shi Han, Yining Chen, Daxin Jiang, and Dongmei Zhang. 2021. Table2Charts: Recommending Charts by Learning Shared Table Representations. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining* (Virtual Event, Singapore) (KDD '21). Association for Computing Machinery, New York, NY, USA, 2389–2399. <https://doi.org/10.1145/3447548.3467279>
- [52] Sujia Zhu, Guodao Sun, Qi Jiang, Meng Zha, and Ronghua Liang. 2020. A survey on automatic infographics and visualization recommendations. *Visual Informatics* 4, 3 (2020), 24–40. <https://doi.org/10.1016/j.visinf.2020.07.002>