

Are We There Yet? A Review on Existing Perceptual Theory and Experiment Support for Visualization Recommendation Systems

Zehua Zeng, Minhui Xie, Matthew Gouzoulis, Leilani Battle

Abstract—A growing body of research focuses on helping users explore complex datasets faster by automatically suggesting visualization designs of possible interest. However, existing visualization recommendation systems only enumerate, rank, and recommend a small group of visualization designs. Our goal is to understand whether there is enough theoretical and experimental knowledge in current literature to inform visualization recommendation systems to assess the entire visualization design space. Thus, in this paper, we present a literature review comparing and ranking the quality of visualization designs in visual perception and human performance. We structure our review by first defining the visualization design space where visualizations must be compared to recommend effective visualization designs. We then perform the review by using a comprehensive schema to record the theoretical and experimental results of visualization comparison, which can also be used to guide the future construction of visualization recommendation systems. To analyze the literature coverage, we develop an interactive tool that can help explore current literature coverage of visualization comparison and identify gaps efficiently and effectively. Based on our findings, we highlight new opportunities and challenges for the community in working towards a comprehensive visualization ranking for informing visualization recommendation systems.

Index Terms—Literature Review, Human Perception, Visualization Recommendation Systems

1 INTRODUCTION

VISUALIZATION design can be a complex process [1], [2]; multiple competing factors must be weighed [3], [4], such as reconciling the size and type of data being visualized [5], [6], [7], the target analysis task being supported [5], [8], and even individual user characteristics [9]. Visualization recommendation algorithms have been developed to reduce user effort in visual analysis and exploration tasks by partially or fully automating the visualization design process (e.g., Show Me [7], Voyager [10], [11], Foresight [12], DeepEye [4], SeeDB [13]). As previous research demonstrates [14], the majority of visualization recommendation algorithms follow the same two-step process of enumerating and then ranking candidate recommendations. We observe that for most of these algorithms, the ranking step relies on either theoretical principles [7], [10], [11] or experimental evaluations to compare visualizations [3]. Specifically, given a pair of visualization designs, visualization recommendation algorithms use heuristics to approximate, compare, and ultimately rank the visual properties represented within these designs. This ranking strategy is then systematically applied to a large search space of possible visualization designs until only a small subset of the top-ranked designs remain, which are then recommended to the user.

However, these algorithms are only as effective as the input data passed to them. Without concrete data or guidelines as inputs, it would be impossible for these algorithms to compare and rank visualization designs to generate high-quality recommendations effectively. A lack of data could cause visualization recommendation

algorithms to be somewhat brittle in practice. These algorithms might resort to recommending only a small set of encoding channels or basic visualization types. For example, Show Me [7] and Voyager systems [10], [11] only recommend charts visualized with position, length, color, area, and shape encodings, while other systems [12], [13], [15] which focus on recommending data queries can only suggest a few visualization types, like bar charts, scatterplots, and line charts. In addition, Zeng et al. [14] find concerning performance issues with existing visualization recommendation algorithms, due primarily to a lack of focus on evaluating how these algorithms compare and rank different visualization designs. These issues point to a fundamental question in visualization recommendation research: *do we have sufficient empirical data (or theory-based rules) to make informed recommendations for different visual analysis tasks and datasets to visualize?*

In this paper, we answer this question through a survey of the existing literature on comparing and ranking visualization designs in terms of visual perception and human performance under different visual analysis tasks. We systematically document the visualization designs studied in each paper, as well as other factors that influence how designs are compared, such as input data characteristics. However, we find significant gaps in the literature where many visualization design pairings have yet to be assessed. For example, we observe a clear emphasis on evaluating the position, length, and color encodings but very little coverage of other encoding types in the literature (e.g., angle, orientation, or density). Furthermore, few papers seem to directly compare a range of alternative visualization designs (i.e., sets of encodings). Without “enough” information on how various visualizations perform compared to each other, the visualization community may continue to fall short in designing new visualization recommendation algorithms that provide measurable improvements in user performance.

- Zehua Zeng, Minhui Xie and Matthew Gouzoulis are with University of Maryland, College Park. E-mail: zheng@umd.edu, minhui.ivy.xie@gmail.com, mgouzou1@terpmail.umd.edu
- Leilani Battle is with University of Washington, Seattle. E-mail: leibatt@cs.washington.edu

Manuscript received April 19, 2005; revised August 26, 2015.

In summary, we make the following contributions in this paper:

- We first review a broad range of the literature (67 papers) on visualization comparison and develop a schema to record the theoretical and experimental results of the comparisons made. This schema can inform the design of new visualization recommendation algorithms.
- We contribute an interactive tool (ARE-WE-THERE-YET?) for analyzing the literature coverage. The tool can help the community quickly find related work and identify gaps in empirical and theoretical research on visualization comparison.
- Finally, we contribute new guidelines for the visualization community and suggest potential paths for future research to address observed challenges in visualization comparison and visualization recommendation.

All of our data are available online: https://osf.io/m6sbp/?view_only=f6fcf140a7bb4c4c9a3e1871f6dc7188.

2 RELATED WORK

In this section, we discuss existing work in visualization recommendation systems, how they define the visualization design spaces and also the design principles behind those systems.

2.1 Visualization Recommendation Systems

Existing visualization recommendation systems can be divided into two main categories according to the strategies they used to rank visualization designs: rule-based or machine learning-based [14], [16]. Rule-based systems utilize either existing theoretical principles [10], [11] in visual perception or propose new metrics [12], [13], [15] to rank visualization designs. For example, Wongsuphasawat et al. [10], [11] use Mackinlay’s principles [6] to make recommendations to the user, prioritizing recommendations based on the breadth of data covered within the visualizations. Vartak et al. [13] use an “interestingness” metric based on deviation in the data to identify visualizations of potential interest. Both Kay et al. [15] and Demiralp et al. [12] apply statistical features of the dataset into their systems for guiding exploratory visual analysis.

Machine learning-based systems design and train their learning models based on (often large) visualization design corpora. For example, Hu et al. [16] trained a deep learning model using an extensive input of the most commonly seen Plotly visualizations and recommend visualization designs for new datasets using the trained model. In a similar spirit, Luo et al. [4] implemented a visualization recommendation system by combining deep learning techniques with hand-written rules. On the other hand, Moritz et al. [3] introduced the Draco system, which enables users to generate relevant visualizations by formulating design requirements as rules passed to a constraint solver. One of the Draco applications, Draco-Learn, was implemented with a training model which learns effectiveness criteria from two prior empirical studies [5], [8].

Based on prior surveys of visualization recommendation systems [14], [17], we can see that most visualization recommendation algorithms rely on either theoretical principles or empirical data in comparing visualization performance. With insufficient data to properly inform recommendation algorithms, the related systems may perform poorly in practical scenarios or fail to recommend an extensive range of visualization designs for exploration [14]. Thus, we aim to estimate the visualization community’s progress in collecting this data and identify specific gaps where more results are needed.

2.2 Principles of Visualization Design

Many works investigate how to best design effective visualizations. Theory works such as Bertin’s visual encoding principles [18], and Mackinlay’s APT work [6] have been highly influential in information visualization research. Cleveland & McGill [19] organized the encoding channels put forth by Bertin from least to most effective in terms of what users can perceive from quantitative data, and validated this ranking in part through visual perception studies. Mackinlay [6] later extended the ranking to include ordinal and nominal data in the APT system. Shneiderman [20]’s task taxonomy then broadens Mackinlay’s work by including data types that were not covered in APT, such as multidimensional data, trees, and networks. The design principles proposed by Bertin, Cleveland & McGill, Mackinlay, and Shneiderman inform the structure of our framework, which focuses on organizing visualization quality comparison tasks based on the underlying visual encodings within a design, rather than by high-level visualization types.

Numerous later experiments build on these foundational theoretical works. For example, the experimental results of Cleveland & McGill were replicated and validated by Heer & Bostock [21] through crowdsourcing of visual perception experiments. Talbot et al. [22] also designed four follow-up experiments on the perception of bar charts to further explore and explain Cleveland & McGill’s results. Their main goal was to understand how different bar chart designs impact analysis task performance. Kim et al. [5] discuss ways to evaluate the effectiveness of twelve 3-encoding visualization designs for different low-level tasks and dataset characteristics. Saket et al. [8] evaluate the effectiveness of basic visualization types for a specific set of analysis tasks.

We aim to assess the overall impact of existing empirical and theoretical work in terms of coverage in our literature survey. Specifically, we record the visualization designs and tasks investigated in these works and analyze the current literature coverage. Potential gaps within the literature would then represent blind spots, where we have insufficient data to determine whether one visualization design is more effective than another, limiting the ability of visualization recommendation systems to rank and recommend a larger range of visualization designs. How to effectively and efficiently fill up those gaps is also one of our main discussions.

2.3 Visualization Design Spaces

The visualization design space must be clearly defined before visualization recommendation algorithms can enumerate and rank candidate visualization designs. For this reason, the visualization design space also defines the scope of coverage for our review of the visualization comparison literature. Prior work uses different methods for defining the visualization design space, which informs our work. For example, Voyager [10], [11] utilizes Vega-Lite [23] specifications to represent visualizations. Draco [3] defines its design space as a set of hard constraints, and visualizations are included in the space only if they satisfy all constraints. DeepEye [4] defines its visualization design space using data queries and limits the space of supported visualization types to bar charts, pie charts, line charts, and scatter plots.

In this work, we first establish the boundaries of the visualization design space based on previous perceptual studies, then adopt a similar framework proposed by Satyanarayan et al. [23] to record the researched visualization designs in the literature.

3 LITERATURE REVIEW METHODS

Our goal is to assess the community’s progress towards designing robust algorithms that can reason intelligently about various visualization designs across analysis tasks and datasets. Thus, we aim to investigate existing theory and experiment papers that focus on evaluating and comparing visualizations in visual perception and user performance. In this section, we describe our methodology rationale and the method for collecting and filtering papers. Then we present the schema we develop to record the visualization ranking results from relevant theoretical and experimental studies.

3.1 Establishing Design Space Boundaries

It is impractical to derive a single visualization recommendation system to cover all possible uses of visualization. However, it is equally impractical to expect visualization users to learn a completely different system for every conceivable visualization use case. We need to strategically define the visualization design space within which we believe a single recommendation system (and algorithm) can actually be effective. Our design space boundaries are informed by existing literature on (1) visualization design spaces, which formally define the range of visualization designs that could be recommended; and (2) visual perception studies, which can be used to identify a subset of designs that can be fairly compared in terms of user performance. We summarize our findings as four constraints on the visualization design space.

Exclude 3D visualizations. As found in previous work, users often have difficulty in perceiving information within 3D visualizations [24]. Moreover, in many cases, multiple linked 2D views prove to be more effective than a single 3D visualization of the same data [25]. Thus, we exclude 3D visualizations from our design space.

Focus on static visualization designs. Although animations and transitions can improve a user’s perception of an underlying dataset [26], many if not most visualizations are still designed without any animations or transitions. Given an apparent lack of data in the literature evaluating the animation and transition design spaces, we do not include these design elements within our visualization design space. On the other hand, the design space of interactions is still an under-explored area in visualization, and enumeration of this space has only recently become viable [23]. In this case, the lack of data and theoretical principles is already evident and does not require an in-depth literature review. As a result, we exclude interactions from our analysis.

Exclude small multiples. Small multiples are designed to present different slices of a multi-dimensional dataset and consist of multiple visualizations side by side, each with the same set of encodings [27]. However, small multiples require significantly more time to interpret than single visualizations [5] since they spread the viewer’s attention and perception of the underlying data across multiple similar but not identical visualizations. Furthermore, charts can be arranged differently in a small multiples view [28]; thus, small multiples are not necessarily perceived the same way as single visualizations [28], [29], [30]. For these reasons, we exclude small multiples views from our visualization design space.

Limit visualizations to three or fewer visual encodings. According to a preliminary literature review, we find in practice that the vast majority of papers study visualization designs with three encodings or less, with 2-3 encodings being the norm. For example, Voyager [10], [11] rarely generates visualizations with more than 3 encodings. Similarly, Tableau’s Show Me feature [7] switches to

using small multiples when recommending encoding channels to more than 2-3 attributes at a time. Thus we limit our visualization design space to consider only designs with 1-3 encodings.

3.2 Method for Collecting Papers

Here, we discuss the assumptions and process behind our search for relevant literature.

3.2.1 Filtering for Relevant Papers

Given the boundaries of the visualization design space established in Sect. 3.1, we seek to understand how much of this space is actually covered (i.e., evaluated) by the existing literature. Given our focus on visualization recommendation, we use the following filters to guide our paper selection process:

Focus on human perception and task performance. An essential facet of visualization recommendation systems is encoding selection, which directly impacts a user’s ability to perceive the underlying information [14], [31]. Even if a visualization system suggests certain attributes to explore, these findings will be inaccessible to the user if the data is presented incorrectly. Therefore, we focus on results that speak to a user’s ability to perceive different visual encodings, as well as differences in user performance across tasks and visualization designs.

Focus on evaluation with regular displays. Although some existing work has researched the effect of display size on visual perception or task performance [21], [32], [33], and some are building new systems to better support different display sizes [34], [35], [36], the vast majority of existing visualization evaluations are still conducted in regular displays (e.g., computer screens). Thus, we focus on reviewing the literature in visualization evaluation and comparison with regular displays.

Emphasize comparisons of different visualization designs. In order for algorithms to select the most relevant visualization design for a given dataset, they must be able to compare and ultimately rank the effectiveness of different designs [14]. To determine which designs should be preferred by these algorithms, we need experimental results that compare different visualization designs, or theoretical rules and guidelines to prune irrelevant designs. To this end, we include any paper in our review that evaluates the user’s ability to effectively perceive and reason about information encoded using one or more visualization designs included in our design space, and compares or provides guidance on ranking pairs of relevant visualization designs.

3.2.2 Paper Sampling and Collection

To initially find relevant papers for our literature review, we checked all papers published in visualization-related conferences in the last 5 years, performed multiple online searches, and asked for recommendations from colleagues. We used the following keywords in our searches: “encoding”, “perception”, “evaluate”, “effectiveness”. We also reviewed the references for each paper found through colleagues or online searches; any relevant papers were also included in our review. In total, we found 147 candidate papers for our literature review.

We then excluded papers that fall outside the boundaries of the visualization design space described in Sect. 3.1 or fail to pass the filters described in Sect. 3.2.1. For example, we excluded papers that only evaluate animated visualizations, 3D visualizations, or small multiples. In another example, Stasko’s work [37] on constructing a value-driven model to evaluate visualizations is

not within the scope of our framework, since comparison of specific visualization designs is not discussed with respect to visual perception. We also exclude papers that reuse data from previous papers (e.g., [3], [38], [39]), since these works do not contribute new information in terms of coverage of the visualization design space. This filtering step excluded 80 of the 147 candidate papers, leaving 67 papers for our analysis.

3.3 Schema for Recording Perceptual Results

To enable a fine-grained analysis of the visualization design space, for each paper we reviewed, we record the paper title, the category of paper, the covered visualization designs, the perceptual tasks measured, and the theoretical or empirical rankings derived for each task. We also record the paper’s authors and the paper link to help users identify and locate the research paper when using our interactive tool (discussed in Sect. 5). We developed a JSON-based schema to record each of these components:

Category: either *theory* or *experiment*, or *hybrid*. A *hybrid* paper presents theoretical hypotheses and experiments to test (at least some of) the proposed hypotheses.

Covered Designs: a list of all researched visualization designs which can fall into our visualization space boundaries.

Other Designs: a list of all researched visualization designs which can **not** be described with our visualization space definition, such as tables, and parallel coordinates.

Tasks: a list of the theoretically discussed or empirically tested tasks. We utilize Amar et al. [40]’s task taxonomy as a starting point and also add more analytics tasks that fall outside the taxonomy to the list (see Table 1).

Results: ranking lists of the performance of visualization designs towards different analysis tasks. Results are separated into experimental results and theoretical rankings.

3.3.1 Visualization Designs

Inspired by existing visualization grammars [23], [41], we use four components to specify each visualization design observed:

Data Attribute: can be either quantitative, nominal, or ordinal.

Encoding Channel: can be position (X- and/or Y-axis), length, angle, area, texture, volume, density, shape, color saturation, color hue, or orientation (see Fig. 1).

Aggregation: defined by the aggregation type and window

Mark Type: can be either point, line, area-circle, area-rect, area-other, box-plot

When selecting encoding channels for our analysis, we start with the encoding channels discussed in the ranking of perceptual tasks proposed by Cleveland & McGill [19] and later extended by Mackinlay [6]. We remove the *connection* and *containment* channels, because we find hardly any papers that evaluate them. We also find that *orientation* has been discussed frequently in the literature (e.g., [42], [43]), and is similar to the *direction* channel proposed by Cleveland & McGill [19] and the *slope* channel mentioned by Mackinlay [6]; thus we combine them into one *orientation* channel. We also split the *position* channel into *positionX* and *positionY*, since there are 2 directions of position in the 2D Cartesian plane, which could impact a user’s perception of these values, bringing the number of encoding channels to 12. To save space, we use an abbreviation to represent each channel. PX is positionX, PY is positionY, L is length, An is angle, Ar is area, T is texture, V is volume, D is density, S is shape, CS is color saturation, CH is color hue and O is orientation, shown in Fig. 1.

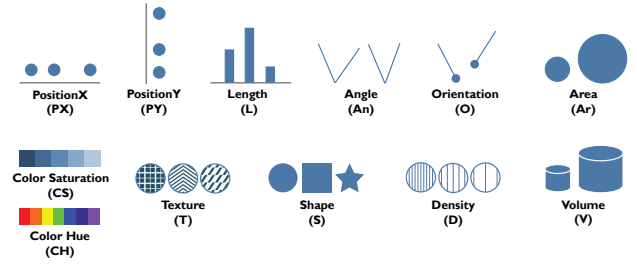


Fig. 1: All encoding channels observed in our literature review.

Since some designs are generated by overlaying multiple visualizations on top of each other, we also incorporate layers in our schema [41]. For example, Albers et al. used a composite graph, a bar chart overlaid on a line chart (as demonstrated in Listing 1), to evaluate supporting aggregate time-series comparison tasks [44]. Each layer is a single visualization design, defined by an encoding list and mark type. Each encoding list consists of the attribute(s) used in the layer, as well as the assigned encoding channel(s) and aggregation details. The encoding channel is defined as a list since multiple encoding channels can be assigned to visualize one attribute in some situations, while aggregation information is stored as an object with the aggregation type and aggregation window. We assign an ID for each recorded visualization design, where “C” means it belongs to covered designs, and “O” means other designs, while “E” means the design is empirically evaluated, “T” means it is theoretically discussed. For example, the ID “CE-4” in Listing 1 means this visualization design is the fourth empirically tested design in the paper [44].

Note that when determining the encoding list for a visualization design, we consider all encodings a user or participant perceives within the design rather than the subset of encodings highlighted by a particular experiment or design rule. For example, when participants are asked to judge whether test marks are using the same or different colors in scatterplots [45], they perceive three encodings, even if only one encoding (color) is permuted.

Listing 1: Example of a covered visualization design, where a bar chart overlaid on a line chart.

```
"CE-4": {
  "layers": [
    {
      "encodings": [
        {
          "attribute-type": "quantitative",
          "aggregation": {
            "aggregation-type": "none",
            "aggregation-window": "none"
          },
          "channels": ["positionY"],
          {
            "attribute-type": "ordinal",
            "aggregation": {
              "aggregation-type": "none",
              "aggregation-window": "none"
            },
            "channels": ["positionX"]
          },
          "marks": ["line"]
        },
        {
          "attribute-type": "quantitative",
          "aggregation": {
            "aggregation-type": "average",
            "aggregation-window": "month"
          },
          "channels": ["length"],
          {
            "attribute-type": "ordinal",
            "aggregation": {
              "aggregation-type": "none",
              "aggregation-window": "none"
            },
            "channels": ["positionX"]
          },
          "marks": ["area-rect"]
        }
      ]
    }
  ]
}
```

TABLE 1: A taxonomy of visual analysis tasks observed in our literature review.

Tasks	Descriptions	Relevant Work
Amar et al.'s [40] Task Taxonomy	1. Retrieve value	Identify values of the specified attributes [5], [8], [43], [46], [47], [48], [49], [50], [51], [52], [53], [54]
	2. Filter	Find data points satisfying the specified conditions [8]
	3. Compute Derived Value	Compute the aggregate value of the specified attributes [8], [51], [55], [56], [57], [58], [59], [60]
	4. Find Extremum Find Maximum Find Minimum	Find data points with an extreme (maximum/minimum) value of the specified attribute [5], [8], [42], [43], [44], [54], [59], [61], [62], [63], [64]
	5. Sort Compare Objects	Rank a set of data points by the specified ordinal metric Compare data points to determine if they are the same (or different) [5], [8], [43], [45], [47], [51], [57], [58], [65], [66], [67]
	6. Determine Range	Find the span of values of the specified attributes [8], [44], [62]
	7. Characterize Distribution	Identify the distribution of given attributes [8], [44], [51], [54], [68]
	8. Find Anomalies	Identify anomalies within the dataset [8], [44], [49], [51], [54], [55], [69], [70], [71]
	9. Cluster Detect Number of Clusters	Find clusters of similar attribute values Detect the number of groups of similar data attribute values [8], [49], [51], [54], [69], [70], [71], [72], [73]
	10. Correlate Estimate Correlation	Determine whether correlation exists within the specified attributes Estimate the correlation coefficient between two given attributes [8], [51], [54], [69], [70], [74], [75], [76], [77], [78], [79], [80], [81], [82]
Added Tasks	11. Compare Derived Values	Compare the aggregate values of multiple attributes [5], [44], [51], [55], [61], [62], [63], [83], [84], [85], [86]
	12. Estimate Trend	Estimate the trend of a given attribute within a specified time [42], [47], [49], [54], [55], [58], [61], [63], [87], [88]
	13. Estimate Difference	Compare and estimate the percentage of difference between the given attributes [19], [21], [22], [52], [53], [57], [89], [90]
	14. Detect Presence	Determine whether the target is present or absent [91]
	15. Locate	Locate the given target [43], [51], [92], [93]
	16. Recognize	Identify the referent from the given representation [51], [54], [55], [62]
	17. Explore Adjacency	Explore the properties of adjacency objects [51], [58]
	18. Judge Similarity	Judge the similarity among multiple visualizations or find the visualization most like the target chart [54], [88], [94], [95], [96], [97], [98]
	19. Compare Chart Structure	Find the visualization is most structured/clustered/clumsy/complex among a group of visualizations [72], [88]

3.3.2 Tasks

Since different experiments and theories focus on different visual analysis tasks, we developed a standardized taxonomy to categorize the many observed tasks. We used a mix of a priori and inductive coding to develop a comprehensive list of visual analysis tasks from the literature. We started with the low-level analysis tasks proposed by Amar et al. as a priori codes [40] and added new task codes to fill in observed gaps, resulting in 19 total task codes (10 a priori codes, 9 inductive codes). Table 1 shows all of the visual analytics tasks we observed in our literature review, as well as their descriptions and the related work that mentions them.

3.3.3 Results

Results are separated into experiment results and theory rankings. They are grouped by the observed metric [99] used for ranking: *accuracy*, *JND* (just noticeable difference), *time* and *user-preference* for experimental ranking and *effectiveness* for theoretical ranking. Under each metric, the visualization designs are sorted based on their performance in supporting each task (from the best to the worst). For visualizations that perform about the same, we group them into a sub-list. While the *rank* list only reflects the ranking among visualization designs, it does not show whether there is a significant difference between two visualization designs in terms of each metric. Thus, we also record a *significance* list that stores pairs of visualizations whose performances are significantly different (the front one performs significantly better than the latter).

We use a hybrid paper [19] as example, as shown in Listing 2. If a task is conducted multiple times, we add an index after the task name to differentiate different task rounds. For example, from Listing 2 we can see that two *estimate-difference* tasks were

conducted. Five designs were involved in the first task, while only two are evaluated in the second task. The *rank* lists show that visualization “CE-1” performed the best in both tasks, and the *significance* lists reveal that there was a significant performance difference between “CE-1” and some other visualization designs.

Listing 2: Example of the result of a hybrid paper [19].

```

"Results": {
  "Experimental": {
    "accuracy": {
      "estimate-difference-1": {
        "rank": ["CE-1", "CE-2", "CE-3", "CE-4", "CE-5"],
        "significance": [["CE-1", "CE-3"], ["CE-1", "CE-4"],
          ["CE-1", "CE-5"], ["CE-2", "CE-4"], ["CE-2",
            "CE-5"], ["CE-3", "CE-4"], ["CE-3", "CE-5"]]],
      "estimate-difference-2": {
        "rank": ["CE-1", "CE-6"],
        "significance": [["CE-1", "CE-6"]]}},
    "Theoretical": {
      "effectiveness": {
        "overall": {
          "ranking": ["CT-1", "CT-2", ["CT-3", "CT-4", "CT-5"],
            "CT-6", "CT-7", ["CT-8", "CT-9"]],
          "significance": []}}}}

```

4 LITERATURE COVERAGE

In this section, we review existing visualization theories and experiments that could be used to compare and rank pairs of visualization designs.

4.1 Theoretical Principles

Theory principles are critical for helping recommendation algorithms prune the design space. For example, theoretical hypotheses,

TABLE 2: Coverage for 2-encoding visualization designs. The papers in italics are *theoretical*, the underlined ones are *hybrid* and the rest are experimental. To save space, we remove all empty rows and columns. Note that there are actually $\sum_{i=1}^{11} = 66$ cells in the full table.

	PX	PY	L	Ar
PY	<i>[51], [54], [81], [100], [101]</i> , [8], [21], [44], [46], [47], [56], [57], [60], [61], [62], [74], [75], [76], [77], [78], [79], [80], [83], [87], [95], [19], [69], [71], [72], [82], [88], [96], [97], [98]			
L	<i>[54]</i> , [8], [21], [22], [43], [44], [46], [47], [48], [52], [53], [56], [57], [60], [62], [66], [19], [65]	[46], [57]		
Ar	[21], [42], [46], [48], [52], [53], [57], [58], [87], [89], [95]	[21], [46], [57], [58], [89]	[59]	
CH	[42], [50], [83], [95]	–	[90], [65]	[8], [46], [75], [90]
CS	[42], [44], [50]	–	[59], [68]	[59], [68]
O	[42]	–	–	[21], [43], [19]
T, S	[42]	–	–	–
D	[48]	–	–	–
V	[46]	[46]	–	–

such as abstract rules, can predict what visualization designs may not be effective and thus may be eliminated early. Informed by theoretical rules, recommendation algorithms can potentially reduce the number of visualization designs that must be compared and ranked during the recommendation process.

1-encoding ranking. Since we are using the perceptual rankings proposed by Cleveland & McGill [19] and Mackinlay [6] to partially define the visualization design space, all of the 1-encoding rankings are covered by their work. To investigate further, we also consider Mackinlay’s expressiveness and effectiveness principles. As Mackinlay suggests, texture and shape encodings are a poor fit for quantitative attributes, and shape is irrelevant to the ordinal data type. Mackinlay argues that size (area and length encodings) and color saturation encodings may be perceived as ordered, so they should not be used for conveying nominal data. Furthermore, Mackinlay also states that only parts of the color spectrum can be perceived as ordinal since the full-color spectrum is not ordered. Thus, these seven attribute-encoding pairs can be considered ineffective and inexpressive: (Q, T/S), (O, S/CH), (N, Ar/L/CS). By applying these perceptual hypotheses, recommendation algorithms could reduce numerous visualization designs that need to be enumerated and ranked.

In addition, we also note an extensive body of theoretical work on evaluating color hue and saturation [65], [91], [102], [103], [104]. For example, Lin et al. [65] introduced an algorithm for automatic selection of semantically-resonant colors to represent data. Fang et al. [104] presented an algorithmic approach for maximizing the perceptual distances among a set of given colors. Bujack et al. [103] surveyed the literature to analyze what features of color are necessary or sufficient to imply perceptual order.

As we observe that 1-encoding visualization designs are rare and cannot deliver much information, the vast majority of the visualization designs include 2-3 encodings. Although existing theory provides a full 1-encoding ranking, we are not sure how exactly these principles can be applied to designs with two or more encodings. For example, according to Mackinlay’s perceptual ranking, *position* and *length* are the top two encodings for visualizing quantitative data, and *position* performs better than *length*. However, it is unclear whether the chart visualizing two quantitative attributes with both position encodings (PX & PY) is more effective than the one visualizing the same attributes with one position (PX/PY) and one length encoding (L).

2-encoding visualization designs. We find that existing theoretical principles provide inadequate coverage of 2-encoding visualization designs; they focus primarily on PX+PY designs, with only one observation of PX/PY+L designs, as shown in Table 2.

Among theories about PX+PY visualization designs, we find that most focus on scatterplots [51], [54], [69], [71], [72], [81], [96], [97], [98], [100], [101] (11 papers total), only three discuss line charts [54], [88], [98] and one bar charts [54].

For example, consider the theory work on *Scagnostics* by Tukey and Tukey [105] and later Wilkinson et al. [101], which is widely used to evaluate the perceptual effectiveness of scatterplots. Though comprehensive for scatterplots, no other types of visualizations are compared. Similarly, Sedlmair and Aupetit [70] proposed a data-driven framework to evaluate how well a measure would predict human judgments specifically for scatterplot perceptual tasks. Sarikaya and Gleicher [51] generated a framework to help designers determine which scatterplot designs are appropriate considering different analysis scenarios. Ryan et al. [88] created a new measure for visual complexity, Pixel Approximate Entropy, for line charts.

3-encoding visualization designs. We only found 5 theory papers that provide concrete guidance for ranking 3-encoding visualization designs; all seem to emphasize color [51], [54], [64], [70], [73]. Szafr et al. [54] reviewed existing literature and provided examples of the ensemble of four common visual encodings (position, area, orientation, and color) on solving four visual aggregation tasks (summary, identification, pattern recognition, and segmentation). Sedlmair and Aupetit [70], on the other hand, proposed a framework to evaluate how well existing measures can predict human judgment of class separability for color-coded scatterplots (PX+PY+CH). Wang et al. [73] focused on optimizing color assignment with respect to the perception of multi-class scatterplots (PX+PY+CH). Cheng et al. [64] design a data-driven color assignment method that can be applied to heat maps, choropleth maps, and diagrams (PX+PY+CH).

Existing theory can help us evaluate the quality of a visualization design in advance of running experiments, but there is a limitation that should not be ignored. Current visualization theory is often unable to provide quantitative estimates for comparing the effectiveness of different visual encodings. For example, Mackinlay suggested that a position encoding should be ranked higher than a length encoding for conveying quantitative data. However, whether

the ranking is significant and by how much is still unclear. Thus, we need empirical data to refine and reinforce the high-level guidelines put forth in existing theory work.

4.2 Experimental Evaluation

Experiments often lead to the formulation of new theories and the refinement of existing ones. Thus experiments are integral to the process of evaluating visualization designs. In this section, we analyze how existing visual perception experiments cover the visualization design space. We summarize observed coverage, then expand on our results by discussing topics that seem to be emphasized (or not) in the literature.

We observed two kinds of experimental evaluation in our literature review: evaluating different design decisions within a single chart type (within-encoding) or comparing multiple designs visualized using different encodings (between-encoding).

4.2.1 Evaluating different design decisions.

Even when the combination of $\{data\ attributes, data\ transformations, encoding\ channels, mark\ type\}$ is fixed, there are still other nuanced design decisions which might affect human perception. We found many empirical work which evaluate whether changing the visualization design slightly may affect task performance [22], [47], [50], [55], [61], [65], [66], [69], [71], [72], [73], [74], [77], [78], [79], [79], [80], [82], [86], [88], [89], [92], [96], [97] (24 papers total). For example, Burlinson et al. [55] conducted a series of experiments to see whether using different combinations of shapes (e.g., stars–open shape, triangles–closed shape) to represent classes in multi-class scatterplots (PX+PY+S) would affect human perception. Schloss et al. [86] investigated whether people’s prediction of quantitative values would be influenced by changing visual features in colormap visualization (PX+PY+CH), such as background colors and color scales. Zhao et al. [66] tested how much the perception of a bar changes when bars are ordered differently in a bar chart. However, these experiment results only provide a more profound but not broader understanding of visualization comparison as guidelines to inform visualization recommendation systems.

4.2.2 Comparing different encodings.

2-encoding visualization designs. In total, there are 19 experimental papers that compare multiple 2-encoding visualization designs. We ground this analysis of comparison coverage using the seminal experiments of Cleveland and McGill [19], which have clearly influenced the choice and trajectory of later empirical works in visualization design evaluation. Cleveland and McGill not only hypothesized an ordering of visual encodings on the basis of graphical perception but also tested parts of this theory through experiments. They used bar charts (PX+L) to assess *position* and *length* encodings and pie charts (Ar+O) for *area* encoding¹. Heer and Bostock [21] replicated these experiments but also adjusted the experiments to make results between length and area encodings comparable.

On the other hand, Saket et al. [8] conducted a crowdsourced experiment to evaluate the effectiveness of five 2-encoding visualization designs across ten analysis tasks: line chart (PX+PY),

bar chart (PX+L), scatterplot (PX+PY), pie chart (Ar+CH) and table (not within our defined visualization space). To evaluate whether using radial charts to visualize daily patterns is effective, Waldner [43] compared radial charts (Ar+O) with bar charts (PX+L) across a mix of low-level and high-level tasks; the results suggest that bar charts are more effective, although the clock-like radial chart is commonly used in our daily life. Chung et al. [42] conducted two crowdsourced empirical studies that focus on the perceptual evaluation of order-ability for a large range of encoding channels, using several 2-encoding designs such as (PX+Ar), (PX+T), (PX+S), etc. They found that specific visual channels are perceived as more ordered (e.g., color saturation) than others (e.g., color hue).

3-encoding visualization designs. We discovered 14 experiment papers that evaluate multiple 3-encoding visualizations designs [5], [45], [46], [49], [55], [62], [63], [64], [67], [73], [75], [84], [85], [86]. For example, Kim and Heer [5] conducted a crowdsourced experiment comparing subject performance across twelve 3-encoding visualization designs, as well as different task types and data distributions; positionX, positionY, area, color saturation, and color hue encodings are compared. Szafir [45] conducted a series of studies that focus on measuring human perceptions of color difference for three common mark types: points, bars, and lines; multi-class scatterplots (PX+PY+CH), bar charts (PX+L+CH), and line charts (PX+PY+CH) are compared. Gramazio [92] conducted experiments with color matrices (PX+PY+CH) and multi-class scatterplots (PX+PY+CH) to learn how the grouping, quantity, and size of visual marks affect human perception and search time.

Ranking both 2- and 3-encoding designs. We also found 3 papers that performed evaluations among mixed 2- and 3-encoding visualization designs [46], [62], [75]. For example, Harrison et al. [75] conducted a crowdsourced experiment to investigate whether Weber’s law could be used to model the perception of correlation in nine commonly used visualization types: scatterplot (PX+PY), area chart (Ar+CH), colored line chart (PX+PY+CH), parallel coordinates, etc. Their results indicated that Weber’s law can model all tested visualizations, but the effectiveness varies according to the underlying data correlation.

4.3 Well-covered Visual Encodings

Here, we discuss encodings currently emphasized in the literature.

4.3.1 Position

Position is the most commonly discussed encoding channel in both the theoretical and experimental literature. Although the PX+PY design could refer to scatterplots or line charts, scatterplots are mentioned significantly more than line charts.

Multi-class scatterplots. Gleicher et al. [84] investigated relative mean value judgments within multi-class scatterplots where color hue, luminance, and shape encodings, or a combination thereof, were used to represent different classes. Wang et al. [73] proposed an approach for assigning colors to multi-class scatterplots based on a set of given colors with the goal of optimizing the perception of scatterplots. Demiralp et al. [94]’s experiment results on estimating perceptual distance for different encoding channels have also provided relevant guidelines for selecting visual encodings for multi-class scatterplots.

Predicting scatterplot perception. Sedlmair and Aupetit [70] constructed a framework that learns how a quality measure could predict human judgments on class separability in 2D scatterplots.

1. Although Cleveland and McGill [19] considered pie charts as relying on an angle encoding channel, more recent research [106] suggests that humans perceive pie charts primarily with area cues. Thus, pie charts are classified as area encodings throughout our categorization.

To predict the perceptual distances between scatterplots as scored by human annotators, Jo and Seo [100] compute 32 representative features to capture the characteristics of bivariable data distributions and pass them as an input to a neural network. In a related vein, ScatterNet [107] exploits deep neural networks to extract semantic features of scatterplot images for similarity calculation.

4.3.2 Length

Visualization designs that combine position and length are also commonly discussed in the experimental literature. We find bar charts to be the overwhelming majority of these designs.

Measuring perceptual bias with bar charts. Godau et al. [56] tested whether there is a bias in the central tendency perceived in bar charts, and they found that the mean value was systematically underestimated in bar charts (but not in scatterplots). Their other experiments also confirmed that the underestimation of the average was not affected by including outliers. Xiong et al. [60] conducted three empirical studies to investigate position perception bias with visualizations containing a single bar/line, multiple bars/lines, and one line with one set of bars. In contrast to the results of Godau et al., they found that the perceived average position was significantly biased in both single line charts and single bar charts. Line positions were consistently underestimated, while bar positions were overestimated. In the experiments involving multiple data series (multiple lines and/or bars), they also observed an effect of “perceptual pull”, where the average position estimate for each series was “pulled” toward the other.

Visualization embellishment. We find a number of works that evaluate the effectiveness of bar charts when styled using methods outside of standard encoding channels [46], [48], [52], [53]. Borgo et al. [46] reported an empirical study on using visual embellishments (e.g., adding informational icons or changing styling on marks) in different visualization designs, including bar charts, line charts, bubble charts, and treemaps. Their results suggest that visual embellishments can help people better remember information but have a negative impact on visual search. Haroz et al. [48], tested how pictographic bar charts impact memory, search performance, and engagement. They found no user costs and even some benefits to using pictographic bar charts; however, superfluous images can still distract users. Skau et al. [52] focused on the effects of visual embellishments on data communication in bar charts. They summarized common embellishments made to bar charts, such as rounded top bar charts, triangle bar charts, capped bar charts, etc. Their results suggest that bar chart embellishments do have an impact on how well data can be communicated within the visualization: even small changes like rounded tops led to higher error rates in perceiving the underlying data.

4.3.3 Color

Color saturation and hue are also well represented in the literature.

Evaluating colormaps. Reda et al. [63] investigated the effects of colormap characteristics and spatial frequency (a proxy for data complexity) on the perception of continuous colormaps. Their results indicate that spatial frequency does impact human judgment of encoded quantities and structures. Reda and Papka [85] evaluate the effectiveness of different colormaps at depicting gradient magnitudes. To evaluate and compare different single-hue and multi-hue colormaps, Liu and Heer [50] conducted crowdsourced experiments testing subject performance on judging relative distances perceived within color triplets. After surveying existing guidelines for

TABLE 3: All possible combinations of data types based on the number of attributes. **Q** means quantitative, **N** means nominal, and **O** means ordinal attribute type. Subscripts are used to distinguish different attributes with the same data type.

# of Attributes (k)	$k = 1$	$k = 2$	$k = 3$
Attribute Combinations			$(Q_1, Q_2, Q_3),$ $(N_1, N_2, N_3),$ $(O_1, O_2, O_3),$
		$(Q_1, Q_2),$ $(N_1, N_2),$ $(O_1, O_2),$	$(Q_1, Q_2, N),$ $(Q, N_1, N_2),$ $(Q_1, Q_2, O),$
	$(Q),$ $(N),$ (O)	$(Q, N),$ $(Q, O),$ (N, O)	$(Q, O_1, O_2),$ $(N_1, N_2, O),$ (N, O_1, O_2) (Q, N, O)

colormap design, Bujack et al. [102] proposed a mathematical framework to describe and assess colormap properties.

Color assignment. There appear to be two existing research directions in the categorical color assignment. One is to assign different colors to data values. Fang et al. [104] developed an algorithm for maximizing the perceptual distances among a given set of categorical data, while Wang et al. [73] focus more on assigning colors to optimize the perception of clusters within multiclass scatterplots. Another research direction in the color assignment is mapping colors to meanings. Lin et al. [65] developed an automatic selection algorithm that assigns semantically relative colors to encode given data. Schloss et al. [86] extend the concept of inferred color mapping from previous work to investigate how the inferred color assignment would be affected by background color.

4.4 Calculating Literature Coverage

To calculate the literature coverage of the design space, we first need to know how many possible visualization designs are within the space. If we only consider the combination of encoding channels, ignoring attribute types, there would be total ${}^{12}C_1 + {}^{12}C_2 + {}^{12}C_3 = 298$ possible combinations of 1-, 2- and 3-encoding designs. However, we only observe 40 different encoding designs mentioned in the literature. As shown in Table 5, theory papers cover 17 (5.7%) out of 298 designs, while experiment papers cover 32 (10.7%), and hybrid papers cover 15 (5%).

On the other hand, if we consider both attribute types and encoding channels, there will be many more variants. First, we only consider the unique combinations of attribute types that can be selected per encoding group, shown in Table 3. We then calculate how many possible unique combinations of encoding channels can be assigned to each combination of attribute types. Though attribute types can be reused across encodings, the encoding channels cannot. For example, when selecting one encoding channel, there are 12 possible options, one per available channel. When selecting 2 encoding channels, there are ${}^{12}P_2 = 132$ possible combinations. Thus, considering both attribute types and encoding channels, there are total 14,028 different visualization designs in the space (also as shown in Table 4):

$${}^3C_1 \times {}^{12}P_1 + {}^4C_2 \times {}^{12}P_2 = 792 + {}^5C_3 \times {}^{12}P_3 = 14,028$$

Based on our findings, although the theory literature covers all 1-encoding designs, few 2-encoding and 3-encoding visualization designs are covered by existing theory (0.38% and 0.03% of

TABLE 4: The number of unique designs within the visualization design space. (Original) is without any filtering, and (APT Rules) is applying Mackinlay’s effectiveness and expressiveness principles, while (Potential) excludes uninteresting combinations of attribute types $\{(N_1, N_2), (O_1, O_2), (N_1, N_2, N_3), (O_1, O_2, O_3)\}$ and (Potential + APT) is applying both filter criteria.

# of attributes	1	2	3	Total
Original	36	792	13,200	14,028
APT Rules	29	509	6,695	7,233
Potential	36	528	10,560	11,124
Potential + APT	29	347	5,471	5,847

all relevant designs, respectively). Furthermore, the experimental literature appears to be imbalanced in assessing different visual encodings. Although position, length, and color encodings are relatively well studied, other encodings are under-explored. We find that 19.4% of 1-encoding and 4.2% of 2-encoding visualization designs are covered by current experiment papers. We find that 14 experiment papers evaluate a subset of 3-encoding designs, but only a tiny percentage of this space is covered (0.16%).

In summary, we find that neither the theory nor experimental literature that we reviewed provide sufficient coverage for visualization recommendation algorithms. As a result, we believe that the robustness and scope of these algorithms will remain limited until more theoretical rules and empirical findings are contributed to fill existing gaps.

5 INTERACTIVE TOOL FOR COVERAGE ANALYSIS

We develop an interactive tool to help analyze the literature coverage. This tool can not only help the community find related

work, but also identify gaps in theoretical and experimental research on visualization comparison.

5.1 Interface Design

Fig. 2 shows the interface of the interactive tool for analyzing the literature coverage, which contains a *filter panel* (top), a *design coverage panel* (left), a *summary view* (middle) and a *detail view* (right). We provide both the source code and a demonstration video of the tool in the OSF repository.

5.1.1 Filter Panel (A)

The filter panel allows users to filter papers by the paper category, tasks, the number of encodings, and also encoding designs. To filter papers by the paper category, tasks, and the number of encodings, users only need to check/uncheck the desired checkboxes, while to filter by encoding designs, users need to select the combination of encodings and then submit the combination. By clicking the (✕) button along with the encoding design, the selection of the specific encoding design would be canceled (as shown in Fig. 2 (A)).

5.1.2 Design Coverage Panel (B)

The design coverage panel shows all of the encoding designs mentioned by the papers satisfying the current filter condition.

5.1.3 Summary View (C)

The summary view displays all papers which satisfy the current filter condition. Each paper view shows the paper category, the paper title, the paper’s authors, and researched tasks.

5.1.4 Detail View (D)

By clicking a paper in the summary view, the detail view will display the detailed information of the specific paper, including covered designs and the theoretical and/or experimental ranking of the covered designs concluded from the studies of the paper.

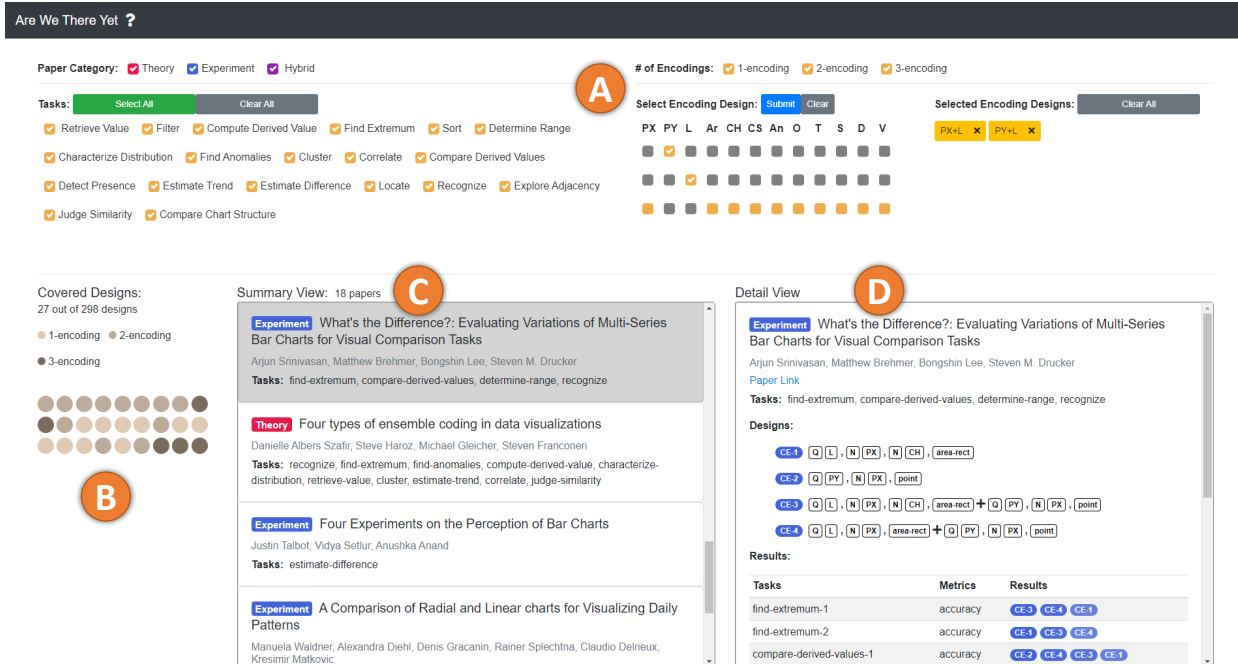


Fig. 2: Interactive tool for analyzing literature coverage. The **filter panel** (A) allows users to filter research papers by paper categories, tasks, numbers of encodings, and encoding designs. The **summary view** (C) displays all of the papers which satisfy the current filter condition and the **design coverage panel** (B) shows all of the visualization designs mentioned by the papers from (C). The **detail view** (D) demonstrates the detailed information of the selected paper.

TABLE 5: The number of encoding designs covered by each paper category.

	1-encoding	2-encoding	3-encoding	Total
Theory	12	2	3	17
Experiment	5	19	8	32
Hybrid	10	4	1	15

5.2 Use Cases

To demonstrate the value of our tool, we present two motivating scenarios of direct benefit to the visualization community: (1) searching for related work in visualization comparison and (2) identifying gaps in literature coverage.

5.2.1 Searching for Related Work

Search by encoding design. Suppose a researcher is searching related work on whether different types of bar charts (vertical, horizontal, grouped, stacked) perform the same in terms of human perception and user performance. By selecting the encoding designs of bar charts in our search interface (PX+L & PY+L), the summary view will show all papers that mention (PX+L) and/or (PY+L) encoding designs. As shown in Fig. 2 (C), there are 18 papers in this space. The researcher can also click on an exciting paper to check the detailed information of the paper.

Search by analysis task. Imagine another researcher aims to develop a new task taxonomy for visual analysis and exploration. She could use the tool to search the relevant work of each task, as well as how the task is described in the literature and how it was conducted to compare different visualization designs.

5.2.2 Identifying Gaps

As demonstrated in our tool, only 40 out of 298 encoding designs are studied in our surveyed 67 papers.

Filtering by paper category. To recreate our findings for theory and experiment coverage in the literature in Sect. 4, we can filter papers by paper category. When we apply category filters in our search interface, we can see that theoretical papers only cover 17 designs, among which 12 are 1-encoding visualization designs, 2 are 2-encoding designs, and the rest are 3-encoding designs. Likewise, we observe that experimental papers cover 32 visualization designs while hybrid papers cover 15 designs (as shown in Table 5).

Filtering by task type. We can also research how many encoding designs are studied under each analysis task by filtering by task type. In this way, we find that some analysis tasks are relatively well researched; for example, there are more than 15 encoding designs have been evaluated with *compute-derived-value*, *find-extremum*, and *estimate-difference* tasks. However, fewer than 3 encoding designs have been evaluated for other tasks, such as *filter*, *detect-presence*, and *compare-chart-structure*.

6 DISCUSSION

In this paper, we present a thorough literature review to investigate how visualization designs are compared and ranked in existing visualization theory and experimental work. We conclude our findings from the survey with a discussion on how current literature coverage can inform visualization recommendation systems. We also outline the challenges we observed in the existing literature and provide some research directions as guidelines for the community towards developing a complete visualization ranking list based on human perception.

6.1 Informing Visualization Recommendation Systems

Although the current literature does not cover the entire visualization design space, we can still derive helpful knowledge to inform visualization recommendation systems.

6.1.1 Initial Pruning

Visualization recommendation systems generally follow a similar recommendation process [14], [108]: first **enumerate** the visualization space to find candidates, then **rank** the candidates based on specific criteria. For example, the perceptual ranking proposed by Mackinlay [6] could be used to identify portions of the visualization design space that are ineffective or inexpressive. By pruning the visualization space based on Mackinlay’s hypotheses, i.e., removing visualizations containing any sub-optimal attribute-encoding pairs (Q, T/S), (O, S/CH), (N, Ar/L/CS), we can reduce the number of visualization designs that need to be **ranked** by a visualization recommendation algorithm. In fact, one of the Draco [3] applications, Draco-APT, utilizes Mackinlay’s effectiveness and expressiveness criteria to exclude visualizations that do not satisfy the rules from its design space.

6.1.2 Ranking Different Encodings

In terms of empirical support, both Saket et al. [8] and Kim et al. [5] conducted experiments to evaluate a larger group of visualization designs. Saket et al. evaluated five common-seen 2-encoding visualizations (scatter, bar, line, pie, table) under all ten analysis low-level task scenarios categorized by Amar et al. [40]. All three data attribute types were tested in the evaluation. Visualizations are compared based on the effectiveness (*accuracy* and *time*), as well as based on the *user-preference* metric. On the other hand, Kim et al. [5] evaluated ten different 3-encoding visualizations involving the same 1 and 2 quantitative nominal attributes, and including *PX*, *PY*, *Ar*, *CS*, and *CH* encodings. Four out of ten task conditions from Amar et al.’s taxonomy are tested. Although the visualization design space is not entirely covered yet, at least for the tested visualization designs, we do have an understanding of which encodings are better for specific tasks based on the empirical results. If we combine data and findings from multiple experiments to cover more of the visualization design space, then we can potentially develop more robust visualization recommendation algorithms, e.g., by translating these findings into a consolidated set of hard constraints in Draco [3].

6.1.3 Refining Visualization Designs

As we mentioned in 4.2, some experiments focus on evaluating different design decisions, like changing the visual features (e.g., background colors, color scales) for colormap visualizations, changing the display order of bars in bar charts, etc. The results derived from this type of experiment can inform visualization recommendation systems of nuanced design decisions, which can be added as hand-tuned rules for finalizing visualization designs. For example, the results from Zhao et al. [66]’s experiments suggest that the height of its neighbors influences the perception of a bar. In addition, the task performance also depends on the number of data points in the visualization and other data characteristics of the dataset. By learning from such knowledge, a recommendation algorithm can manipulate the order and number of bars in the bar chart towards better task performance.

6.2 Challenges towards Comprehensive Ranking

6.2.1 Gaps in Visualization Comparison Coverage

As demonstrated in Section 4, there exist numerous visualization designs in the space that remained unexplored.

Gaps in theoretical work. All twelve encodings are ranked (or considered ineffectively) in theory work based on how well they can convey different types of data (quantitative, nominal, ordinal) according to Mackinlay [6]’s perceptual hypotheses. However, it is unclear how different encoding channels work together to deliver more information in multi-encoding visualizations. Although some theory work has researched multi-encoding visualizations, they mainly focus on scatterplots (PX+PY), like predicting human judgments in some specific tasks with scatterplots. That is to say, different encodings are rarely theoretically compared in multi-encoding visualization designs.

Gaps in experimental work. Experimental evaluation, on the other hand, covers more multi-encoding visualization designs compared to theoretical work. However, existing evaluations mainly focus on half of the encoding channels (*positionX*, *positionY*, *length*, *area*, *color-hue*, *color-saturation*) proposed in Mackinlay [6]’s work. The other half of the encodings (like *texture*, *volume*, *shape*, *angle*, *density*) are hardly evaluated (by 1 or 2 papers maximum) or never tested.

6.2.2 Inconsistencies and Conflicts in the Literature

We find a number of inconsistencies and conflicts in the literature, which may hinder our ability to compare and reuse existing results for visualization recommendation systems.

Between theories and experiments. We observe inconsistencies in the way visual encodings are ascribed to specific designs. As previously mentioned, Cleveland and McGill treated pie charts as primarily angle encodings [19]; however, more recent experimental work suggests that pie charts are perceived more as area encodings [106]. Furthermore, although theory work can help us potentially prune the visualization design space, the theoretical hypotheses might not be necessarily “correct” in the practical environment. Seven attribute-encoding pairs ((Q, T/S), (O, S/CH), (N, Ar/L/CS)) are considered ineffective and inexpressive based on Mackinlay [6]’s work, however, a more recent evaluation [42] shows different results. Mackinlay suggests that *texture* and *shape* are not relevant to quantitative data type, but according to the results from Chung et al.’s experiments, both *texture* and *shape* performed better than *orientation* conveying quantitative data in *estimate-trend* and *find-extremum* tasks.

Between different experiments. Even when experiments were similar, we may find contradictory results. Even though Godau et al. [56] and Xiong et al. [60] both conducted experiments to test human bias in perceiving average position for length (bar charts) and position encodings (scatterplot or line charts), they had completely different results. Godau et al. only found underestimation in bar charts but no bias for point positions (scatterplots). However, Xiong et al. found significant bias in both bar charts and line charts, where line positions were underestimated while bar positions were overestimated. In another example, Harrison et al. [75] find Weber’s law to be a convincing model for how people perceive data correlations, however in a re-analysis of the same data, Kay and Heer [39] find Weber’s Law not to be a good fit. It is natural in science to improve upon existing results and theories; however, there is currently no easy way to identify and track these discrepancies within the literature and translate

them into concrete improvements to visualization recommendation systems.

6.3 Potential Research Directions

Here, we propose some potential research paths that the visualization community can focus on to better cover the visualization design space and enhance the visualization recommendation process.

6.3.1 Paths for Theoretical Work

More theoretical hypotheses for pruning the design space.

Theory papers are critical for helping us prune the design space. For example, by eliminating visualization designs containing any sub-optimal attribute-encoding pairs (Q, T/S), (O, S/CH), (N, Ar/L/CS) proposed by Mackinlay [6], the number of visualization designs (considering both attribute types and encoding channels) is almost down by half (as shown in Table 4). If more theoretical researches like Mackinlay’s [6] are proposed to exclude the ineffective visualization designs, it will help improve the efficiency of visualization recommendation systems since much fewer visualization designs need to be enumerated and ranked during the recommendation process. For example, the combinations of some attribute types (like (N_1, N_2) , (O_1, O_2) , (N_1, N_2, N_3) , (O_1, O_2, O_3)) might not be interesting to visualize. If we can exclude these combinations of attributes from the design space, the space will be much smaller, and also it will become more efficient to compare visualization designs within it (as shown in Table 4 at the last two rows).

More theoretical hypotheses for comparing multi-encoding visualization designs.

As we mentioned that, according to our observation, existing visualization theory only provides substantial coverage for our 1-encoding visualization designs; it covers little of our 2-encoding and almost none of our 3-encoding visualization designs. However, only with the 1-encoding ranking, it is hard to know how different encodings are with or against each other and whether the performance of the combination of multiple encodings still ranks in the same order. Thus, it would be beneficial for the community, especially for informing visualization recommendation algorithms, to propose a perceptual ranking of multi-encoding designs.

Refining theoretical hypotheses based on recent experimental work.

Although influential theory work, such as Bertin [18]’s, Cleveland & McGill [19]’s and Mackinlay [6]’s principles, is still serving essential roles in informing and guiding current visualization recommendation systems [3], [10], [11], these theories were proposed over 30 years ago. On the other hand, more recent experimental work [5], [8], [42] shows slightly different visualization comparison results compared to the theoretical hypotheses. Refining the essential theory work based on recent experimental results would greatly help improve the performance of visualization recommendation systems.

6.3.2 Paths for Experimental Work

More comparative experiments on different encodings and visualization designs.

Experiments are critical for building and refining our understanding of how humans perceive different visualization designs. However, the experimental literature only emphasizes half of the twelve encodings, leaving the rest of the encodings under-explored. These significant gaps in the literature might suggest a lack of rigorous understanding of how these visualization designs perform in the real world. Thus, evaluating and comparing the performance of uncared-for encodings (*texture*,

volume, shape, orientation, angle, density) under different analysis task scenarios would contribute more “ground truth” evidence to adequately ground our attempts at automating the visualization design process.

Re-visiting controversial results. As we discussed in the previous section 6.2.2, there exist inconsistencies and conflicts between theories and experiments, as well as between different experiments. These existing conflicts among theory and practice require more nuanced consideration about the design process as well as the evaluation metrics. Redesigning experiments to test out visualizations where controversial results exist with a more comprehensive comparison between nuanced design decisions and involving more metrics can provide more precise ideas on ranking and finalizing visualizations under different circumstances.

Uplift the Role of Replication Studies in Visualization Research. Given the discrepancies we have observed, we argue that the findings of both visualization theory and experimental research should be treated as hypotheses only until subsequent replication experiments converge on a consistent set of results. Furthermore, we argue that replication experiments should be held in high regard within the visualization community regardless of whether their findings reinforce or challenge our current assumptions, since either way, they enable us to verify our understanding of how people perceive and use visualizations. We need this validation to ensure that later visualization recommendation algorithms are built upon a solid foundation of theoretical and empirical visualization research; we should reward these types of studies accordingly within our community.

7 LIMITATIONS & FUTURE WORK

Informed by existing work on visualization design space and visual perception studies, we excluded (1) 3D visualizations, (2) visualizations with animations or interactions, and (3) small multiples from our defined visualization space. However, it would be interesting to expand our work to include some (or all) of these excluded designs back to the space by the time they are well researched in the literature. With a more extensive design space where the comparison and ranking of visualization designs are well studied, one can design a robust visualization recommendation system to output effective visualization designs under different analytic scenarios.

Given our initial goal is to understand how different visualization designs (specially different encodings) are ranked in the existing theoretical and experimental work, our schema only records {*attribute types, aggregation details, encoding channels, mark type*} for each covered design (details in Listing 1). For other more detailed design decisions, we only support adding a note to specify. For example, Liu and Heer [50] evaluated and compared how subjects perform relative distance judgments among different color schemes (single- and multi-hue). Since the difference among these visualizations (PX+CS for single-hue, PX+CH for multi-hue) is the use of different color scales, we add a note under each design to specify which color scale was used. During the review process, we have added various notes to specify not only the color scales [50], [63], [85] but also whether visualization embellishments exist [48], [52], [53], and whether the two encodings for one attribute is conflicting or redundant [84]. However, it is hard to parse all these notes automatically. Thus, one potential avenue for future work could be improving the schema by adding more design specifications to specify each visualization

design in more detail. This improvement can also help compare different nuanced design decisions.

We also note that by focusing on visual perception, we are unable to account for other factors that may influence the overall effectiveness of a visualization design, such as visual aesthetics [109], intuition and metaphors [110], as well as user background and preferences [9]. Developing a broader framework encompassing both visual perception and these other factors would be exciting future work.

Our literature review contributes a detailed record of how different visualization designs are compared and ranked in the current theoretical and experimental work. This record not only specifies all researched visualization designs, but also keeps track of the ranking of their performance (*accuracy, JND, time, user-preference*) under different task scenarios. A next step to extend this work could be applying the findings to develop a better visualization recommendation system, such as adopting the recommendation strategy based on the user’s current analysis task or analysis intent.

ACKNOWLEDGMENTS

Authors would like to thank the BAD Lab and our paper reviewers for their thoughtful feedback. This work was supported in part by NSF award IIS-1850115 and an Adobe Research Award.

REFERENCES

- [1] M. Chen and H. Leitte, “An information-theoretic framework for visualization,” *IEEE Transactions on Visualization and Computer Graphics*, vol. 16, pp. 1206–1215, 2010.
- [2] M. Chen and A. Golan, “What may visualization processes optimize?” *IEEE Transactions on Visualization and Computer Graphics*, vol. 22, 06 2015.
- [3] D. Moritz, C. Wang, G. L. Nelson, H. Lin, A. M. Smith, B. Howe, and J. Heer, “Formalizing visualization design knowledge as constraints: Actionable and extensible models in draco,” *IEEE Transactions on Visualization and Computer Graphics*, vol. 25, no. 1, pp. 438–448, Jan 2019.
- [4] Y. Luo, X. Qin, N. Tang, and G. Li, “Deepeye: Towards automatic data visualization,” in *2018 IEEE 34th International Conference on Data Engineering (ICDE)*, April 2018, pp. 101–112.
- [5] Y. Kim and J. Heer, “Assessing effects of task and data distribution on the effectiveness of visual encodings,” in *Computer Graphics Forum*, vol. 37, no. 3. Wiley Online Library, 2018, pp. 157–167.
- [6] J. Mackinlay, “Automating the design of graphical presentations of relational information,” *ACM Trans. Graph.*, vol. 5, no. 2, p. 110–141, Apr. 1986. [Online]. Available: <https://doi.org/10.1145/22949.22950>
- [7] J. Mackinlay, P. Hanrahan, and C. Stolte, “Show me: Automatic presentation for visual analysis,” *IEEE Transactions on Visualization and Computer Graphics*, vol. 13, no. 6, pp. 1137–1144, Nov 2007.
- [8] B. Saket, A. Ender, and Ç. Demiralp, “Task-based effectiveness of basic visualizations,” *IEEE Transactions on Visualization and Computer Graphics*, vol. 25, no. 7, pp. 2505–2512, July 2019.
- [9] C. Ziemkiewicz, A. Ottley, R. Crouser, K. Chauncey, S. Su, and R. Chang, “Understanding visualization by understanding individual users,” *Computer Graphics and Applications, IEEE*, vol. 32, 11 2012.
- [10] K. Wongsuphasawat, D. Moritz, A. Anand, J. Mackinlay, B. Howe, and J. Heer, “Voyager: Exploratory analysis via faceted browsing of visualization recommendations,” *IEEE Transactions on Visualization and Computer Graphics*, vol. 22, no. 1, pp. 649–658, Jan 2016.
- [11] K. Wongsuphasawat, Z. Qu, D. Moritz, R. Chang, F. Ouk, A. Anand, J. Mackinlay, B. Howe, and J. Heer, *Voyager 2: Augmenting Visual Analysis with Partial View Specifications*. New York, NY, USA: Association for Computing Machinery, 2017, p. 2648–2659. [Online]. Available: <https://doi.org/10.1145/3025453.3025768>
- [12] Ç. Demiralp, P. J. Haas, S. Parthasarathy, and T. Pedapati, “Foresight: Recommending visual insights,” *CoRR*, vol. abs/1707.03877, 2017. [Online]. Available: <http://arxiv.org/abs/1707.03877>
- [13] M. Vartak, S. Rahman, S. Madden, A. Parameswaran, and N. Polyzotis, “Seedb: Efficient data-driven visualization recommendations to support visual analytics,” *Proc. VLDB Endow.*, vol. 8, no. 13, p. 2182–2193, Sep. 2015. [Online]. Available: <https://doi.org/10.14778/2831360.2831371>

- [14] Z. Zeng, P. Moh, F. Du, J. Hoffswell, T. Y. Lee, S. Malik, E. Koh, and L. Battle, "An evaluation-focused framework for visualization recommendation algorithms," *IEEE Transactions on Visualization and Computer Graphics (IEEE VIS)*, vol. to appear, no. X, p. X–X, October 2021.
- [15] A. Key, B. Howe, D. Perry, and C. Aragon, "Vizdeck: Self-organizing dashboards for visual analytics," in *Proceedings of the 2012 ACM SIGMOD International Conference on Management of Data*, ser. SIGMOD '12. New York, NY, USA: Association for Computing Machinery, 2012, p. 681–684. [Online]. Available: <https://doi.org/10.1145/2213836.2213931>
- [16] K. Z. Hu, M. A. Bakker, S. Li, T. Kraska, and C. A. Hidalgo, "Vizml: A machine learning approach to visualization recommendation," *CoRR*, vol. abs/1808.04819, 2018. [Online]. Available: <http://arxiv.org/abs/1808.04819>
- [17] S. Zhu, G. Sun, Q. Jiang, M. Zha, and R. Liang, "A survey on automatic infographics and visualization recommendations," *Visual Informatics*, vol. 4, no. 3, pp. 24–40, 2020.
- [18] J. Bertin, *Semiology of Graphics*. University of Wisconsin Press, 1983.
- [19] W. S. Cleveland and R. McGill, "Graphical perception: Theory, experimentation, and application to the development of graphical methods," *Journal of the American Statistical Association*, vol. 79, no. 387, pp. 531–554, 1984.
- [20] B. Shneiderman, "The eyes have it: a task by data type taxonomy for information visualizations," in *Proceedings 1996 IEEE Symposium on Visual Languages*, Sep. 1996, pp. 336–343.
- [21] J. Heer and M. Bostock, "Crowdsourcing graphical perception: Using mechanical turk to assess visualization design," in *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, ser. CHI '10. New York, NY, USA: Association for Computing Machinery, 2010, p. 203–212. [Online]. Available: <https://doi.org/10.1145/1753326.1753357>
- [22] J. Talbot, V. Setlur, and A. Anand, "Four experiments on the perception of bar charts," *IEEE Transactions on Visualization and Computer Graphics*, vol. 20, no. 12, pp. 2152–2160, Dec 2014.
- [23] A. Satyanarayan, D. Moritz, K. Wongsuphasawat, and J. Heer, "Vega-lite: A grammar of interactive graphics," *IEEE Transactions on Visualization and Computer Graphics*, vol. 23, no. 1, pp. 341–350, Jan 2017.
- [24] M. Tory, A. E. Kirkpatrick, M. S. Atkins, and T. Moller, "Visualization task performance with 2d, 3d, and combination displays," *IEEE Transactions on Visualization and Computer Graphics*, vol. 12, no. 1, pp. 2–13, Jan. 2006. [Online]. Available: <https://doi.org/10.1109/TVCG.2006.17>
- [25] J. J. Van Wijk and E. R. Van Selow, "Cluster and calendar based visualization of time series data," in *Proceedings 1999 IEEE Symposium on Information Visualization (InfoVis '99)*, Oct 1999, pp. 4–9.
- [26] J. Heer and G. Robertson, "Animated transitions in statistical data graphics," *IEEE Transactions on Visualization and Computer Graphics*, vol. 13, no. 6, pp. 1240–1247, Nov 2007. [Online]. Available: <https://doi.org/10.1109/TVCG.2007.70539>
- [27] E. R. Tufte and P. R. Graves-Morris, *The visual display of quantitative information*. Graphics press Cheshire, CT, 1983, vol. 2, no. 9.
- [28] B. Ondov, N. Jardine, N. Elmqvist, and S. Franconeri, "Face to face: Evaluating visual comparison," *IEEE Transactions on Visualization & Computer Graphics*, 2019. [Online]. Available: <http://www.umi.acs.umd.edu/~elm/projects/face2face/face2face.pdf>, [PDFhttps://doi.org/10.1109/TVCG.2018.2864884](https://doi.org/10.1109/TVCG.2018.2864884), DOI
- [29] B. Ondov, F. Yang, M. Kay, N. Elmqvist, and S. Franconeri, "Revealing perceptual proxies with adversarial examples," *IEEE Transactions on Visualization & Computer Graphics*, vol. 28, no. 1, 2021. [Online]. Available: <http://users.umi.acs.umd.edu/~elm/projects/perceptual-proxies/revealing-proxies.pdf>, [PDFhttps://osf.io/2re7b/](https://osf.io/2re7b/), OSF(materials)
- [30] N. Jardine, B. Ondov, N. Elmqvist, and S. Franconeri, "The perceptual proxies of visual comparison," *IEEE Transactions on Visualization & Computer Graphics*, vol. 26, no. 1, 2020. [Online]. Available: <http://users.umi.acs.umd.edu/~elm/projects/perceptual-proxies/perceptual-proxies.pdf>, PDF
- [31] Y. Kim, K. Wongsuphasawat, J. Hullman, and J. Heer, *GraphScape: A Model for Automated Reasoning about Visualization Similarity and Sequencing*. New York, NY, USA: Association for Computing Machinery, 2017, p. 2628–2638. [Online]. Available: <https://doi.org/10.1145/3025453.3025866>
- [32] C. Andrews, A. Endert, and C. North, *Space to Think: Large High-Resolution Displays for Sensemaking*. New York, NY, USA: Association for Computing Machinery, 2010, p. 55–64. [Online]. Available: <https://doi.org/10.1145/1753326.1753336>
- [33] A. Dasgupta and R. Kosara, "Pargnostics: Screen-space metrics for parallel coordinates," *IEEE Transactions on Visualization and Computer Graphics*, vol. 16, no. 6, pp. 1017–1026, 2010.
- [34] M. Brehmer, B. Lee, P. Isenberg, and E. K. Choe, "Visualizing ranges over time on mobile phones: A task-based crowdsourced evaluation," *IEEE Transactions on Visualization and Computer Graphics*, vol. 25, no. 1, pp. 619–629, 2019.
- [35] T. Horak, S. K. Badam, N. Elmqvist, and R. Dachsel, *When David Meets Goliath: Combining Smartwatches with a Large Vertical Display for Visual Data Exploration*. New York, NY, USA: Association for Computing Machinery, 2018, p. 1–13. [Online]. Available: <https://doi.org/10.1145/3173574.3173593>
- [36] S. K. Badam, F. Amini, N. Elmqvist, and P. Irani, "Supporting visual exploration for multiple users in large display environments," in *2016 IEEE Conference on Visual Analytics Science and Technology (VAST)*, 2016, pp. 1–10.
- [37] J. Stasko, "Value-driven evaluation of visualizations," in *Proceedings of the Fifth Workshop on Beyond Time and Errors: Novel Evaluation Methods for Visualization*, ser. BELIV '14. New York, NY, USA: Association for Computing Machinery, 2014, p. 46–53. [Online]. Available: <https://doi.org/10.1145/2669557.2669579>
- [38] X. Fu, Y. Wang, H. Dong, W. Cui, and H. Zhang, "Visualization assessment: A machine learning approach," in *2019 IEEE Visualization Conference (VIS)*. IEEE, 2019, pp. 126–130.
- [39] M. Kay and J. Heer, "Beyond weber's law: A second look at ranking visualizations of correlation," *IEEE Transactions on Visualization and Computer Graphics*, vol. 22, no. 1, pp. 469–478, 2015.
- [40] R. Amar, J. Eagan, and J. Stasko, "Low-level components of analytic activity in information visualization," in *IEEE Symposium on Information Visualization, 2005. INFOVIS 2005*. IEEE, 2005, pp. 111–117.
- [41] H. Wickham, "A layered grammar of graphics," *Journal of Computational and Graphical Statistics*, vol. 19, no. 1, p. 3–28, 2010.
- [42] D. H. Chung, D. Archambault, R. Borgo, D. J. Edwards, R. S. Laramée, and M. Chen, "How ordered is it? on the perceptual orderability of visual channels," in *Computer Graphics Forum*, vol. 35, no. 3. Wiley Online Library, 2016, pp. 131–140.
- [43] M. Waldner, A. Diehl, D. Gracanin, R. Splechtna, C. Delrieux, and K. Matkovic, "A comparison of radial and linear charts for visualizing daily patterns," *IEEE Transactions on Visualization and Computer Graphics*, vol. 26, pp. 1033–1042, 2019.
- [44] D. Albers, M. Correll, and M. Gleicher, "Task-driven evaluation of aggregation in time series visualization," in *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, ser. CHI '14. New York, NY, USA: Association for Computing Machinery, 2014, p. 551–560. [Online]. Available: <https://doi.org/10.1145/2556288.2557200>
- [45] D. A. Szafir, "Modeling color difference for visualization design," *IEEE Transactions on Visualization and Computer Graphics*, vol. 24, no. 1, pp. 392–401, Jan 2018.
- [46] R. Borgo, A. Abdul-Rahman, F. Mohamed, P. W. Grant, I. Reppa, L. Floridi, and M. Chen, "An empirical study on using visual embellishments in visualization," *IEEE Transactions on Visualization and Computer Graphics*, vol. 18, no. 12, pp. 2759–2768, Dec 2012.
- [47] M. Correll, E. Bertini, and S. Franconeri, "Truncating the y-axis: Threat or menace?" *CoRR*, vol. abs/1907.02035, 2019. [Online]. Available: <http://arxiv.org/abs/1907.02035>
- [48] S. Haroz, R. Kosara, and S. L. Franconeri, "Isotype visualization: Working memory, performance, and engagement with pictographs," in *Proceedings of the 33rd Annual ACM Conference on Human Factors in Computing Systems*, ser. CHI '15. New York, NY, USA: Association for Computing Machinery, 2015, p. 1191–1200. [Online]. Available: <https://doi.org/10.1145/2702123.2702275>
- [49] R. Kanjanabose, A. Abdul-Rahman, and M. Chen, "A multi-task comparative study on scatter plots and parallel coordinates plots," *Computer Graphics Forum*, 6 2015.
- [50] Y. Liu and J. Heer, *Somewhere Over the Rainbow: An Empirical Assessment of Quantitative Colormaps*. New York, NY, USA: Association for Computing Machinery, 2018, p. 1–12. [Online]. Available: <https://doi.org/10.1145/3173574.3174172>
- [51] A. Sarikaya and M. Gleicher, "Scatterplots: Tasks, data, and designs," *IEEE Transactions on Visualization and Computer Graphics*, vol. PP, pp. 1–1, 08 2017.
- [52] D. Skau, L. Harrison, and R. Kosara, "An evaluation of the impact of visual embellishments in bar charts," *Computer Graphics Forum*, vol. 34, 06 2015.
- [53] D. Skau and R. Kosara, "Readability and Precision in Pictorial Bar Charts," in *EuroVis 2017 - Short Papers*, B. Kozlikova, T. Schreck, and T. Wischgold, Eds. The Eurographics Association, 2017.

- [54] D. Albers Szafir, S. Haroz, M. Gleicher, and S. Franconeri, "Four types of ensemble coding in data visualizations," *Journal of Vision*, vol. 16, p. 11, 03 2016.
- [55] D. Burlinson, K. Subramanian, and P. Goolkasian, "Open vs. closed shapes: New perceptual categories?" *IEEE Transactions on Visualization and Computer Graphics*, vol. 24, no. 1, pp. 574–583, 2018.
- [56] C. Godau, T. Vogelgesang, and R. Gaschler, "Perception of bar graphs - a biased impression?" *Comput. Hum. Behav.*, vol. 59, no. C, p. 67–73, Jun. 2016. [Online]. Available: <https://doi.org/10.1016/j.chb.2016.01.036>
- [57] P. Mylavarapu, A. Yalcin, X. Gregg, and N. Elmqvist, "Ranked-list visualization: A graphical perception study," in *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*, ser. CHI '19. New York, NY, USA: Association for Computing Machinery, 2019. [Online]. Available: <https://doi.org/10.1145/3290605.3300422>
- [58] S. Nusrat, M. Alam, and S. Kobourov, "Evaluating cartogram effectiveness," *IEEE Transactions on Visualization and Computer Graphics*, vol. PP, 04 2015.
- [59] C. Perin, T. Wun, R. Pusch, and S. Carpendale, "Assessing the graphical perception of time and speed on 2d+time trajectories," *IEEE Transactions on Visualization and Computer Graphics*, vol. 24, no. 1, pp. 698–708, Jan 2018.
- [60] C. Xiong, C. Ceja, C. Ludwig, and S. Franconeri, "Biased average position estimates in line and bar graphs: Underestimation, overestimation, and perceptual pull," *IEEE Transactions on Visualization and Computer Graphics*, vol. 26, 08 2019.
- [61] W. Aigner, C. Kainz, R. Ma, and S. Miksch, "Bertin was right: An empirical evaluation of indexing to compare multivariate time-series data using line plots," *Computer Graphics Forum*, vol. 30, pp. 215–228, 03 2011.
- [62] A. Srinivasan, M. Brehmer, B. Lee, and S. M. Drucker, *What's the Difference? Evaluating Variations of Multi-Series Bar Charts for Visual Comparison Tasks*. New York, NY, USA: Association for Computing Machinery, 2018, p. 1–12. [Online]. Available: <https://doi.org/10.1145/3173574.3173878>
- [63] K. Reda, P. Nalawade, and K. Ansah-Koi, "Graphical perception of continuous quantitative maps: The effects of spatial frequency and colormap design," in *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems*, ser. CHI '18. New York, NY, USA: Association for Computing Machinery, 2018. [Online]. Available: <https://doi.org/10.1145/3173574.3173846>
- [64] S. Cheng, W. Xu, W. Zhong, and K. Mueller, "A data-driven approach for mapping multivariate data to color," *IEEE Transactions on Visualization and Computer Graphics*, vol. PP, 08 2016.
- [65] S. Lin, J. Fortuna, C. Kulkarni, M. Stone, and J. Heer, "Selecting semantically-resonant colors for data visualization," in *Proceedings of the 15th Eurographics Conference on Visualization*, ser. EuroVis '13. Chichester, GBR: The Eurographs Association John Wiley Sons, Ltd., 2013, p. 401–410.
- [66] M. Zhao, H. Qu, and M. Sedlmair, "Neighborhood perception in bar charts," in *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*, ser. CHI '19. New York, NY, USA: Association for Computing Machinery, 2019. [Online]. Available: <https://doi.org/10.1145/3290605.3300462>
- [67] S. Smart and D. A. Szafir, "Measuring the separability of shape, size, and color in scatterplots," in *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*, ser. CHI '19. New York, NY, USA: Association for Computing Machinery, 2019. [Online]. Available: <https://doi.org/10.1145/3290605.3300899>
- [68] S. Redmond, "Visual cues in estimation of part-to-whole comparisons," *2019 IEEE Visualization Conference (VIS)*, pp. 1–5, 2019.
- [69] L. Micallef, G. Palmas, A. Oulasvirta, and T. Weinkauf, "Towards perceptual optimization of the visual design of scatterplots," *IEEE Transactions on Visualization and Computer Graphics*, vol. 23, no. 06, pp. 1588–1599, Jun 2017.
- [70] M. Sedlmair and M. Aupetit, "Data-driven evaluation of visual quality measures," in *Computer Graphics Forum*, vol. 34, no. 3. Wiley Online Library, 2015, pp. 201–210.
- [71] Y. Wang, Z. Wang, T. Liu, M. Correll, Z. Cheng, O. Deussen, and M. Sedlmair, "Improving the robustness of scagnostics," *IEEE Transactions on Visualization and Computer Graphics*, vol. 26, no. 1, pp. 759–769, 2020.
- [72] M. Abbas, M. Aupetit, M. Sedlmair, and H. Bensmail, "Clustme: A visual quality measure for ranking monochrome scatterplots based on cluster patterns," *Computer Graphics Forum*, vol. 38, no. 3, pp. 225–236, 1 2019.
- [73] Y. Wang, X. Chen, T. Ge, C. Bao, M. Sedlmair, C. Fu, O. Deussen, and B. Chen, "Optimizing color assignment for perception of class separability in multiclass scatterplots," *IEEE Transactions on Visualization and Computer Graphics*, vol. 25, no. 1, pp. 820–829, Jan 2019.
- [74] M. E. Doherty, R. B. Anderson, A. M. Angott, and D. S. Klopfer, "The perception of scatterplots," *Perception & Psychophysics*, vol. 69, no. 7, pp. 1261–1272, 2007.
- [75] L. Harrison, F. Yang, S. Franconeri, and R. Chang, "Ranking visualizations of correlation using weber's law," *IEEE Transactions on Visualization and Computer Graphics*, vol. 20, no. 12, pp. 1943–1952, Dec 2014.
- [76] J. Li, J.-B. Martens, and J. J. van Wijk, "Judging correlation from scatterplots and parallel coordinate plots," *Information Visualization*, vol. 9, no. 1, p. 13–30, Mar. 2010. [Online]. Available: <https://doi.org/10.1057/ivs.2008.13>
- [77] J. Meyer and D. Shinar, "Estimating correlations from scatterplots," *Human Factors*, vol. 34, no. 3, pp. 335–349, 1992.
- [78] J. Meyer, M. Taieb, and I. Flascher, "Correlation estimates as perceptual judgments," *Journal of Experimental Psychology: Applied*, vol. 3, no. 1, p. 3, 1997.
- [79] V. Sher, K. G. Bemis, I. Llicardi, and M. Chen, "An empirical study on the reliability of perceiving correlation indices using scatterplots," *Comput. Graph. Forum*, vol. 36, no. 3, p. 61–72, Jun. 2017. [Online]. Available: <https://doi.org/10.1111/cgf.13168>
- [80] R. F. Strahan and C. J. Hansen, "Underestimating correlation from scatterplots," *Applied Psychological Measurement*, vol. 2, no. 4, pp. 543–550, 1978.
- [81] F. Yang, L. T. Harrison, R. A. Rensink, S. L. Franconeri, and R. Chang, "Correlation judgment and visualization features: A comparative study," *IEEE Transactions on Visualization and Computer Graphics*, vol. 25, no. 3, pp. 1474–1488, March 2019.
- [82] R. Rensink and G. Baldrige, "The perception of correlation in scatterplots," *Comput. Graph. Forum*, vol. 29, pp. 1203–1210, 06 2010.
- [83] M. Correll, D. Albers, S. Franconeri, and M. Gleicher, "Comparing averages in time series data," in *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, ser. CHI '12. New York, NY, USA: Association for Computing Machinery, 2012, p. 1095–1104. [Online]. Available: <https://doi.org/10.1145/2207676.2208556>
- [84] M. Gleicher, M. Correll, C. Nothelfer, and S. Franconeri, "Perception of average value in multiclass scatterplots," *IEEE Transactions on Visualization and Computer Graphics*, vol. 19, no. 12, pp. 2316–2325, 10 2013.
- [85] K. Reda and M. E. Papka, "Evaluating gradient perception in color-coded scalar fields," in *2019 IEEE Visualization Conference (VIS)*, Oct 2019, pp. 271–275.
- [86] K. B. Schloss, C. Gramazio, A. T. Silverman, M. L. Parker, and A. S. Wang, "Mapping color to meaning in colormap data visualizations," *IEEE Transactions on Visualization and Computer Graphics*, vol. 25, pp. 810–819, 2019.
- [87] M. Correll and J. Heer, "Regression by eye: Estimating trends in bivariate visualizations," in *Proceedings of the 2017 CHI Conference on Human Factors in Computing Systems*, ser. CHI '17. New York, NY, USA: Association for Computing Machinery, 2017, p. 1387–1396. [Online]. Available: <https://doi.org/10.1145/3025453.3025922>
- [88] G. Ryan, A. Mosca, R. Chang, and E. Wu, "At a glance: Pixel approximate entropy as a measure of line chart complexity," *CoRR*, vol. abs/1811.03180, 2018. [Online]. Available: <http://arxiv.org/abs/1811.03180>
- [89] N. Kong, J. Heer, and M. Agrawala, "Perceptual guidelines for creating rectangular treemaps," *IEEE Transactions on Visualization and Computer Graphics*, vol. 16, pp. 990–8, 01 2011.
- [90] R. Kosara, "The impact of distribution and chart type on part-to-whole comparisons," in *21st Eurographics Conference on Visualization, EuroVis 2019 - Short Papers, Porto, Portugal, June 3-7, 2019*, J. Johansson, F. Sadlo, and G. E. Marai, Eds. Eurographics Association, 2019, pp. 7–11. [Online]. Available: <https://doi.org/10.2312/evs.20191162>
- [91] C. G. Healey, "Choosing effective colours for data visualization," in *Proceedings of Seventh Annual IEEE Visualization'96*. IEEE, 1996, pp. 263–270.
- [92] C. C. Gramazio, K. B. Schloss, and D. H. Laidlaw, "The relation between visualization size, grouping, and user performance," *IEEE Transactions on Visualization and Computer Graphics*, vol. 20, no. 12, pp. 1953–1962, 2014.
- [93] D. Haehn, J. Tompkin, and H. Pfister, "Evaluating 'graphical perception' with cnns," *IEEE Transactions on Visualization and Computer Graphics*, vol. 25, no. 1, pp. 641–650, 2019.
- [94] Ç. Demiralp, M. S. Bernstein, and J. Heer, "Learning perceptual kernels for visualization design," *IEEE Transactions on Visualization and Computer Graphics*, vol. 20, no. 12, pp. 1933–1942, Dec 2014.

- [95] A. Gogolou, T. Tsandilas, T. Palpanas, and A. Bezerianos, “Comparing similarity perception in time series visualizations,” *IEEE Transactions on Visualization and Computer Graphics*, vol. 25, no. 1, pp. 523–533, Jan 2019.
- [96] J. Matute, A. Telea, and L. Linsen, “Skeleton-based scagnostics,” *IEEE Transactions on Visualization and Computer Graphics*, vol. 24, pp. 542–552, 2018.
- [97] A. V. Pandey, J. Krause, C. Felix, J. Boy, and E. Bertini, “Towards understanding human similarity perception in the analysis of large sets of scatter plots,” in *Proceedings of the 2016 CHI Conference on Human Factors in Computing Systems*, ser. CHI ’16. New York, NY, USA: Association for Computing Machinery, 2016, p. 3659–3669. [Online]. Available: <https://doi.org/10.1145/2858036.2858155>
- [98] Y. Wang, F. Han, L. Zhu, O. Deussen, and B. Chen, “Line graph or scatter plot? automatic selection of methods for visualizing trends in time series,” *IEEE Transactions on Visualization and Computer Graphics*, vol. 24, no. 02, pp. 1141–1154, feb 2018.
- [99] E. Bertini, A. Tatu, and D. Keim, “Quality metrics in high-dimensional data visualization: An overview and systematization,” *IEEE Transactions on Visualization and Computer Graphics*, vol. 17, no. 12, pp. 2203–2212, 2011.
- [100] J. Jo and J. Seo, “Disentangled representation of data distributions in scatterplots,” in *2019 IEEE Visualization Conference (VIS)*, Oct 2019, pp. 136–140.
- [101] L. Wilkinson, A. Anand, and R. Grossman, “Graph-theoretic scagnostics,” in *IEEE Symposium on Information Visualization, 2005. INFOVIS 2005.*, Oct 2005, pp. 157–164.
- [102] R. Bujack, T. L. Turton, F. Samsel, C. Ware, D. H. Rogers, and J. Ahrens, “The good, the bad, and the ugly: A theoretical framework for the assessment of continuous colormaps,” *IEEE Transactions on Visualization and Computer Graphics*, vol. 24, no. 1, pp. 923–933, 2018.
- [103] R. Bujack, T. L. Turton, D. H. Rogers, and J. P. Ahrens, “Ordering perceptions about perceptual order,” in *2018 IEEE Scientific Visualization Conference (SciVis)*, Oct 2018, pp. 32–36.
- [104] H. Fang, S. Walton, E. Delahaye, J. Harris, D. A. Storchak, and M. Chen, “Categorical colormap optimization with visualization case studies,” *IEEE Transactions on Visualization and Computer Graphics*, vol. 23, no. 1, pp. 871–880, Jan 2017.
- [105] J. W. Tukey and P. A. Tukey, “Computer graphics and exploratory data analysis: An introduction,” *The Collected Works of John W. Tukey: Graphics: 1965-1985*, vol. 5, p. 419, 1988.
- [106] R. Kosara, “Evidence for area as the primary visual cue in pie charts,” in *2019 IEEE Visualization Conference (VIS)*, Oct 2019, pp. 101–105.
- [107] Y. Ma, A. K. H. Tung, W. Wang, X. Gao, Z. Pan, and W. Chen, “Scatternet: A deep subjective similarity model for visual analysis of scatterplots,” *IEEE Transactions on Visualization and Computer Graphics*, vol. 26, no. 3, pp. 1562–1576, 2020.
- [108] K. Wongsuphasawat, D. Moritz, A. Anand, J. Mackinlay, B. Howe, and J. Heer, “Towards a general-purpose query language for visualization recommendation,” in *Proceedings of the Workshop on Human-In-the-Loop Data Analytics*, ser. HILDA ’16. New York, NY, USA: ACM, 2016, pp. 4:1–4:6. [Online]. Available: <http://doi.acm.org/10.1145/2939502.2939506>
- [109] E. R. Tufte, *The Visual Display of Quantitative Information*. Cheshire, Connecticut: Graphics Press, 2001.
- [110] C. Ziemkiewicz and R. Kosara, “The shaping of information by visual metaphors,” *IEEE Transactions on Visualization and Computer Graphics*, vol. 14, no. 6, pp. 1269–1276, 2008.