# Generative AI for Visualization: State of the Art and Future Directions

Yilin Ye[a,b], Jianing Hao[a], Yihan Hou[a], Zhan Wang[a], Shishi Xiao[a], Yuyu Luo[a,b], Wei Zeng[a,b]

[a] *The Hong Kong University of Science and Technology (Guangzhou), Guangzhou, Guangdong, China*
[b] *The Hong Kong University of Science and Technology, Hong Kong SAR, China*

**Abstract**

Generative AI (GenAI) has witnessed remarkable progress in recent years and demonstrated impressive performance in various generation tasks in different domains such as computer vision and computational design. Many researchers have attempted to integrate GenAI into visualization framework, leveraging the superior generative capacity for different operations. Concurrently, recent major breakthroughs in GenAI like diffusion model and large language model have also drastically increase the potential of GenAI4VIS. From a technical perspective, this paper looks back on previous visualization studies leveraging GenAI and discusses the challenges and opportunities for future research. Specifically, we cover the applications of different types of GenAI methods including sequence, tabular, spatial and graph generation techniques for different tasks of visualization which we summarize into four major stages: data enhancement, visual mapping generation, stylization and interaction. For each specific visualization sub-task, we illustrate the typical data and concrete GenAI algorithms, aiming to provide in-depth understanding of the state-of-the-art GenAI4VIS techniques and their limitations. Furthermore, based on the survey, we discuss three major aspects of challenges and research opportunities including evaluation, dataset, and the gap between end-to-end GenAI and generative algorithms. By summarizing different generation algorithms, their current applications and limitations, this paper endeavors to provide useful insights for future GenAI4VIS research.

*Keywords:* Visualization, Generative AI

## 1. Introduction

VizDeck [1]. Visualization is a process of rendering graphical representations of spatial or abstract data to assist exploratory data analysis. Recently, many researchers have attempted to apply artificial intelligence (AI) for visualization tasks [2, 3, 4, 5, 6]. Particularly, as visualization essentially involves representations and interactions for raw data, many visualization researchers have started to adopt the rapidly developing generative AI (GenAI) technology, a type of AI technology that empowers the generation of synthetic content and data by learning from existing man-made samples [7, 8]. GenAI has come to the foreground of artificial intelligence in recent years, with profound and widespread impact on various research and application domains such as artifact and interaction design (*e.g.* [9, 10, 11]).

Recently, multi-modal AI generation model such as Stable Diffusion [12] or DaLL-E 2 [13] enable laymen users without traditional art and design skills to easily produce high-quality digital paintings or designs with simple text prompts. In natural language generation, large language models like GPT [14] and LLaMa [15] also demonstrate astounding power of conversation, reasoning and knowledge embedding. In computer graphics, recent models like DreamFusion [16] also shows impressive potential in 3D generation. GenAI's unique strength lies in its flexible capacity to model data and generate designs based on implicitly embedded knowledge gleaned from real-world data. This characteristic positions GenAI as a transformative force capable of alleviating the workload and complexity associated with traditional computational methods, and extending the diversity of design with more creative generated results than previous methods.

The burgeoning potential of GenAI is particularly evident in its ability to enhance and streamline operations throughout the data visualization process. From data processing to the mapping stage and beyond, GenAI can play a pivotal role in tasks such as data inference and augmentation, automatic visualization generation, and chart question answering. For instance, the automatic visualization generation has been a longstanding research focus predating the current wave of GenAI methods, offering non-expert users an efficient means of conducting data analysis and crafting visual representations (*e.g.*, [17, 18]). Traditionally, automatic visualization approaches relied on expert-designed rules rooted in design principles [19]. However, these methods were shackled by the constraints of knowledge-based systems [20], struggling to comprehensively incorporate expert knowledge within convo-

luted rules or oversimplified objective functions. The advent of GenAI introduces a paradigm shift, promising not only increased efficiency but also a more intuitive and accessible approach to visualization in an era marked by unprecedented technological advancements.

Despite the impressive capability of GenAI, when applied to visualization it can face many challenges because of its unique data structure and analytic requirements. For example, the generation of visualization images is significantly different from generation of natural or artistic images. First, the evaluation of GenAI for visualization tasks is more complex than natural image generation as many factors beyond image similarity need to be considered, such as efficiency [21] and data integrity [22]. Second, compared to general GenAI tasks trained on large datasets with simple annotations, the diversity and complexity of visualization tasks demand more complex training data [23], which is harder to curate. Third, the gap between the traditional visualization pipeline with strong rule-based constraints makes it difficult to fully integrate with end-to-end GenAI methods. These unique characteristics makes it less straightforward to leverage the latest pre-trained GenAI models in general domain to empower visualization-specific generation. Therefore, it is important to understand how previous works have utilized GenAI for various visualization applications, what challenges are met and especially how the GenAI methods are adapted to the tasks.

Although some previous surveys have covered the use of AI in a general sense for visualization [3], to the best of our knowledge, no study has focused on comprehensive review of GenAI methods used in visualization. This survey extensively reviews the literature and summarizes the AI-powered generation methods developed for visualization. We categorize the various GenAI methods according to the concrete tasks they address, which correspond to different stages of visualization generation. In this way, we manage to collect 81 research papers on GenAI4VIS. We particularly focus on the different algorithms used in specific tasks in the hope of helping researchers understand the state-of-the-art technical development as well as challenges. We also discuss and highlight potential research opportunities.

This paper is structured as follows. Section 2 outlines the scope and taxonomy of our survey with definition of key concepts. Starting from Section 3 to Section 6, each section corresponds to a stage in the visualization pipeline where GenAI has been used. Specifically, Section 3 concerns the use of GenAI for data enhancement. Section 4 summarizes works leveraging GenAI for visual mapping generation. Section 5 focuses on how GenAI is uti-

lized for stylization and communication with visualization. Section 6 covers GenAI techniques to support user interaction. Each subsection in Section 3 to Section 6 covers a specific task in the stage. Instead of listing the works one by one, the structure of the subsection is divided into two parts: data & algorithm and discussion, for a comprehensive understanding of how the current GenAI method works for data of certain structures and what remains challenging for GenAI in particular tasks. Finally, Section 7 discusses some dominant challenges and research opportunities for future research.

## 2. Scope and Taxonomy

### 2.1. Scope and Definition

Generative AI (GenAI) is a type of AI technique that generates synthetic artifacts by analyzing training examples; learning their patterns and distribution; and then creating realistic facsimiles. GenAI uses generative modeling and advances in deep learning (DL) to produce diverse content at scale by utilizing existing media such as text, graphics, audio, and video [7, 8]. A key feature of GenAI is that it generates new content by learning from data instead of explicit programs.

**GenAI methods categorization.** Despite the differences between different domain targets of generation ranging from text, code, multi-media to 3D generation, the particular algorithms of generation actually depend on the data structures which show common characteristics across different domains. Particularly, in GenAI4VIS applications, categorization based on data structures can facilitate more concrete understanding of the algorithms in relation to the different types of data involved in different visualization tasks. Here, we provide an overview of different types of GenAI in terms of typical data structures associated with data visualization.

- **Sequence Generation**: This category includes the generation of ordered data, such as text, code, music, videos, and time-series data. Sequence generation models, like LSTMs and Transformers, can be used to create content with a sequential or temporal structure.

- **Tabular Generation**: This category covers the generation of structured data in the form of rows and columns, such as spreadsheets or database tables. Applications include data augmentation, anonymization, and data imputation.
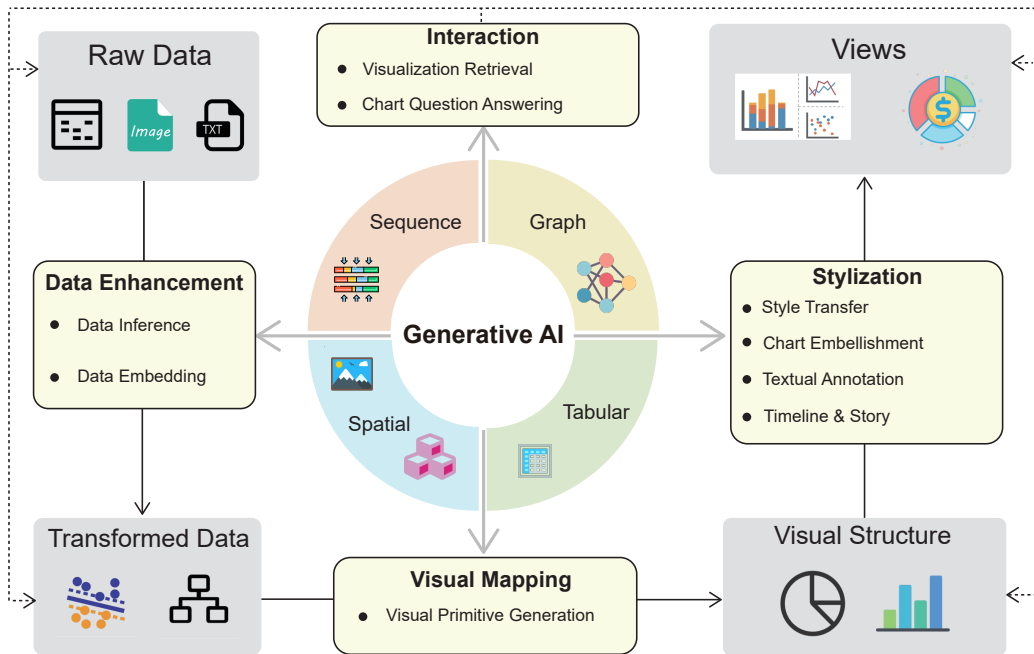
4

Figure 1: The overview of GenAI4VIS applications for different visualization tasks, including data enhancement, visual mapping generation, stylization and interaction tasks.

- **Graph Generation**: This category involves generating graph and network structures, such as social networks, molecular structures, or recommendation systems. Models like Graph Neural Networks (GNNs) and Graph Convolutional Networks (GCNs) can be used to generate or manipulate graph-structured data.

- **Spatial Generation**: This category encompasses the generation of both 2D images and 3D models. These data have the common characteristics of spatial data in 3D or 2D projection in Euclidean spaces, which can be represented as pixels, voxels or points with 2D/3D coordinates. 2D generation includes image synthesis, style transfer, and digital art, while 3D generation covers computer graphics, virtual reality, and 3D printing. Techniques like GANs, VAEs, and PointNet [24] can be used for creating 2D and 3D content.

**GenAI4VIS tasks categorization.** To categorize and organize the collected articles, we are inspired by the classical visualization pipeline describ-

ing different essential stages [25]. However, as GenAI is utilized in broader scenarios different from traditional operations, we also modify the pipeline to encompass some latest research topics. including **data enhancement**, **visual mapping generation**, **stylization**, and **interaction**. Notably, the data transformation part is generalized to the concept of **data enhancement** inspired by the terminology in the study by McNabb et al. [26]. In addition, as few GenAI for visualization works focus on the basic view transformation, we replace this part with a broader concept of **stylization & communication**. Under different stages we further categorize the works into specific tasks, as shown in Figure 1.

- **Data enhancement**. Data enhancement refers to the process of improving the quality or completeness of the data or enhancing the feature representation of the data for subsequent visualization. This can involve data augmentation, embedding or other transformations to make it more suitable for visualization.

- **Visual mapping generation**. This refers to the use of algorithms and software tools to generate visualizations automatically without extensive manual intervention. Automatic visual mapping generation allows users to leverage knowledge about how to create appropriate visualization as common wisdom to reduce the workload and man-made violation of design principles.

- **Stylization**. Extending the concept of presentation in [27], we define stylization in visualization, which involves the application of design principles and aesthetic choices to make the visualization more engaging and effective in conveying information. It includes decisions about color schemes, fonts, layout, and other visual or textual elements to enhance the information-assisted visualization [20].

- **Interaction**. In the context of data visualization, interaction refers to the dynamic engagement and communication between users and the visualized data. It involves the ability of users to manipulate, explore, and interpret visual representations. This can involve various forms of interactivity, such as graphical interactions like zooming, panning, clicking and natural language interaction like chart question answering.

Earlier methods for these tasks focus on rule-based algorithms with complex expert-designed rules reflecting design principles, which is still effective in

6

many applications such as colormap generation [28]. Some studies also leverage optimization-based methods to minimize expert-defined explicit objective functions. However, these types of methods differ from GenAI methods in that they are top-down and do not learn from real-world data. To narrow down the scope of our survey, we exclude all previous generative algorithms that are purely based on rules or optimization.

**Relation between different GenAI methods and tasks.** Due to the wide range of diverse applications in GenAI4VIS, there is no clear-cut one-to-one relation between the type of GenAI methods and the tasks. Nevertheless, we can observe some interesting correlation. First, sequence generation is mostly applied in visual mapping or interaction-related tasks. This is because GenAI such as translation models and the latest LLMs or vision-language model are useful in generating sequence of code specifying visual mapping or sequence of interaction flow and output. Second, tabular generation is mostly used in data enhancement. This is because tabular data with attribute columns are the most common initial input data to visualization, which benefits from data enhancement like surrogate data generation for subsequent tasks. Next, graph generation is also mostly used in data enhancement because data inference and augmentation can facilitate subsequent analysis of graph data. However, despite its relatively rare use, it holds great potential for visual mapping and stylization because graphical structure such as knowledge graph or scene graph can benefit optimization of visual encodings and layout. Finally, spatial generation is mostly applied in data enhancement and stylization tasks. This is because 2D and 3D data such as images and volumetric data are also common types of input for VIS4AI and SciVis applications, while the embellishment of basic charts into stylized charts relies on image-based generation methods. Figure 2 illustrates the relation between GenAI4VIS tasks and methods with a sankey diagram and exemplifies the specific data types that are involved in different methods. Table 1 further list the detailed methodologies for each data structure and task.

*2.2. Related Survey*

Some previous surveys cover the applications of artificial intelligence or machine learning in general to information visualization or scientific visualization [2, 3, 58, 4, 5]. Wu et al. [3] surveyed the development of artificial intelligence technologies applied to information visualization, focusing on three major aspects including visualization data and representation (what), goals
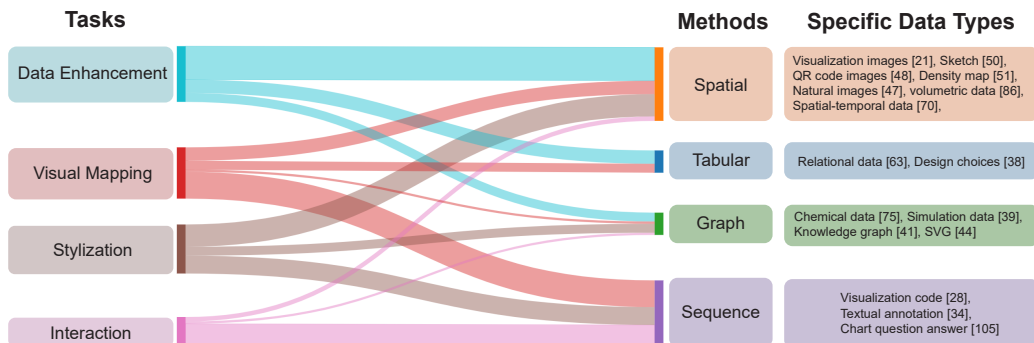
Figure 2: Relation between tasks and methods and examples of specific data types involved in different methods.

(why) and specific tasks (how) which concern the use of AI. Other previous surveys touch upon the use of AI in more specific sub-areas of visualization research, such as natural language interface and data story telling [27, 59, 60]. Shen et al. [27] summarizes all the existing technologies supporting natural language interface to different stages of data visualization, including both traditional rule-based techniques and recent AI-powered methods. Bartolomeo et al. [61] envisions the potential of GenAI applied to different stages of visualization, mostly focused on interviewing experts and discussing usage scenarios. Another closely relevant survey [62] recently focuses on the two-way relationship between visualization and foundation AI models, which include some large scale GenAI models like GPT. In comparison, the focus of our survey is on the concrete technical advancements and challenges of GenAI models for visualization applications.

No previous survey has been dedicated to the various types of GenAI methods in the context of visualization tasks. Our survey aims to provide a specialized and comprehensive overview of the GenAI techniques that have been used to generate various data or intermediate representations that are useful for visualization. We also outline challenges and opportunities for future research on GenAI for VIS.

### 2.3. Survey Methodology

We combine search-based and reference-driven methods to discover relevant literature. We first collect relevant papers from previous surveys and recent works. Then we expand the list by going through the papers' ref-

Table 1: Examples of Specific GenAI4VIS methods applied to different tasks and data types.

| | Data Enhancement | Visual Mapping | Stylization | Interaction |
|---|---|---|---|---|
| **Sequence** | - | RNN [29], Deep Q Network [30], Transformer [31], LLM [32] | LLM [33], RL [34], Detection Network+Template [35] | Detection-based Models [36], Vision-Language Models [37] |
| **Tabular** | Table-GAN [38] | FFN [39], Enumeration+Scoring Network [17] | - | - |
| **Graph** | GNN [40], Latent Traversal [41] | KG embedding [42] | Graph Latent [43], Graph Style Extraction Network [44] | Graph Contrastive Learning [45] |
| **Spatial** | VAE [46], GAN [47], DRL [48], BASNet [49], ISN [50] | Faster R-CNN [51], GAN [52] | Color Extraction Network [53], Siamese Network [54], Diffusion [22], RL [55] | Triplet Autoencoder [56], Contrastive Learning [57] |

erences and citations. In this process, we also supplement the results by searching with key words in the titles of previously collected papers in both the ACM and IEEE libraries as well as arxiv. In the paper selection process, we manually filter out the non-GenAI traditional methods such as purely rule-based or optimization-based methods, stressing the key characteristics of GenAI which has learned from real data in self-supervised pre-training or supervised training stages. In total we collect 81 papers as listed in Table 2 utilizing GenAI for visualization tasks. As shown in Table 2, different types of GenAI4VIS techniques include sequence generation, tabular generation, spatial generation and graph generation, which can benefit visualization tasks in different stages such as data enhancement, visual mapping, stylization and interaction. Sequence and spatial generation are more often used as they concern the general visualization code, images and natural language interaction. We acknowledge that our search method may not be exhaustive due to the manually collection through keyword search and citation traversal. Therefore, this survey mainly provides a comprehensive overview of state-of-the-art GenAI4VIS methods, where application papers with similar methods may not be enumerated.

Table 2: Survey taxonomy and example papers. We classify the collected GenAI4VIS papers into different sub-tasks along visualization pipeline and further classify different GenAI methods into sequence, tabular, spatial and graph generation.

| Tasks | Subtasks | Description | Examples |
|---|---|---|---|
| Data Enhance | Data Inference | Increase samples or dimensions | tabular [63] [64] [38] [65] [66] [67], spatial [46] [68] [48] [69] [70] [71] [47] [72] [73] [69] [48] [74], graph [75] [41] [76] [40] |
| | Data Embedding | Embed data to hide information | spatial [50] [49] [77] |
| Visual Mapping | Visual Primitive | Generate basic visual structures | sequence [29] [30] [78] [31] [79] [80] [32] [81] [82] [83] [84] [85], tabular [39] [17] [86] [18], spatial [51] [52] [87] [88] [89] [90], graph [42] |
| Stylization | Style Transfer | Imitate styles of examples | spatial [53] [91] [92] [93] [54] [94], graph [43] [95] [44] |
| | Embellishment | Generate infographics | spatial [32] [22] [96] |
| | Text Annotation | Add information with text | sequence [35, 33, 97] |
| | Timeline & Story | Generate data story | sequence [98] [99] [100] [101] [34], spatial [55], graph [100] |
| Interact | Retrieval | Find similar charts | spatial [57] [56], graph [45] |
| | CQA | Answer questions about chart | sequence [102] [103], [104] [36] [105] [106] [107] [108] [37] [109] |

## 3. Data Enhancement

### 3.1. Data Inference

GenAI can be useful for inferring unobserved data items or data features based on the distributions and feature values of existing data, such as data interpolation, data augmentation and super-resolution, which we use a general term data inference to describe, as these tasks all aim at inferring unseen data.

### 3.1.1. Graph Generation

**Data.** GenAI for graph data inference is commonly applied in the domain of chemical data.

- *Chemical data.* GenAI-powered data interpolation has been used for interactive exploration of chemical data [110, 75] to assist discovery of new molecule structures. For example, ChemoVerse [75] is an interactive system that leverages interpolation powered by GenAI to help

10

experts understand AI drug design models and verify potential new designs.

- *Graph simulation data.* Graph generation can also be applied to inference of certain physical simulation data which can be modeled as graph data structure, such as ocean simulation data [40].

*Method.* Typical methods include GNN and latent space traversal:

- *Graph Neural Network (GNN).* GNN has been developed to model data that can be represented as graphs [111]. By extracting graph features with operations like graph convolution, GNN can be applied to a wide range of non-Euclidean data with complex relationships, which can also benefit some visualization tasks such as structure-aware visualization retrieval [45] and data reconstruction [40]. For example, GNN-Surrogate [40] is proposed to reconstruct ocean simulation data based on simulation parameters for efficient parameter space exploration. Particularly, because training an end-to-end model that directly reconstruct the full high-resolution ocean simulation data is expensive, the authors proposed to construct an intermediate graph representation for adaptive resolution. The hierarchical graphs are constructed with a series of operations including edge-weighted graph construction, graph hierarchy generation and hierarchical tree cutting. GNN-Surrogate, which is an up-sampling graph generator, first transforms input parameters into latent vector. Then the latent vector is reshaped into initial graph, which is passed through multiple steps of graph convolutions with residual connections. Specifically, graph convolution generates the features of each node by weighted sum of features in the previous layer for the node and all its neighbors. In chemical data interpolation, sometimes special latent space traversal algorithm need to be developed to generate desirable intermediate samples, because the direct linear interpolation assumes that the latent space is flat and Euclidean [112], which may poorly model the complex structure, particularly for chemical molecular data. To address this challenge, some researchers develop special traversal methods [75, 41, 76]. For example, ChemoVerse [75] introduces a manifold traversal algorithm. To find a path going through regions of interest, a k-d tree is built based on the Jacobian distances of all points of interest and additional user-specified

11

constraints. Then $A^*$ algorithm is combined with Yen's algorithm [113] to find the shortest path in the k-d tree. Subsequently, data interpolation is performed along this path by sampling points at equal interval and decoding the latent vectors into full-fledged molecular structures with the generation model.

*3.1.2. Tabular Generation*

**Data**. GenAI can be used to synthesize surrogate data for subsequent tasks, such as privacy protection. Specifically, to protect users' data, oftentimes many institutions would not reveal the real data to the public, causing lack of data for domain-specific analysis tasks. Instead, some studies aim at generating surrogate data similar to real data which can be freely used to test downstream tasks such as visualization and query [63, 64, 38, 65, 66, 67]. The data are typically relational data.

- *Relational data.* Relational data is the most basic form of data for visualization which are often stored in tabular format comprising data items in rows and multi-dimensional attributes in columns. Surrogate data generation studies mainly focus on tabular relational data.

**Method**. A common method for tabular data generation is GAN.

- *Generative Adversarial Network (GAN).* For example, in recent years, some researchers attempt to generate relational data similar to real data with GANs [63, 38, 65, 66]. The architecture of GANs consists of a generator and a discriminator. The adversarial training scheme where the generator progressively learn to generate more realistic data that can deceive the discriminator enables GANs to model the distribution of real data. For example, table-GAN [38] builds upon the basic deep convolutional GAN (DCGAN) [114] framework and tailor the generation to tabular data. Specifically, first the tabular records are converted into square matrix to accommodate convolution operation. In addition to the original generator and discriminator, table-GAN also incorporates a classifier network which learns the correlation between categorical labels and other attributes from the table. This serves to maintain the consistency of values in the generated records. Moreover, besides the original adversarial loss, the authors design a new information loss, which measures the first order and second order statistical

difference between the high-dimensional embedding vectors before the sigmoid function in the discriminator. However, purely GAN-based methods still leak important features of user data as they are directly trained on real data. To address this risk, SERD [65] seeks to generate similar data while preserving the key privacy information in real data. Specifically, SERD manages to satisfy the differential privacy guarantee conditions by using fake entities satisfying the same vectorized similarity constraint of entities in real datasets.

### 3.1.3. Spatial Generation

*Data*. GenAI can be used to infer spatial data such as imagery, volumetric or spatial-temporal data.

- *Imagery data.* Some studies apply GenAI to image data inference, including emoji images [46], medical images [115], natural images [116, 117], etc. The inferred data are used for subsequent tasks like visual interpolation [46, 118], super-resolution [115], 3D reconstruction [68], object detection [48] and semantic segmentation [69].

- *Volumetric data.* Volumetric data is a type of data that represents information in three-dimensional space, which is widely used in various fields such as biology, geology, and physics. Some studies use GenAI for volumetric data super-resolution to address the problem of low data quality [70]. Others works focus on volumetric data reconstruction through generation model. For example, DeepOrganNet [68] applies GenAI to reconstructing and visualizing high-fidelity 3D organ models based on input of merely single-view medical images.

- *Spatial-temporal data.* Spatial-temporal data such as flow data is a type of data that combines both spatial and temporal components. It involves information that varies not only in space (location) but also over time. GenAI can be applied to data extrapolation of spatial-temporal data. For example, Wiewel et al. [71] leverages GenAI to model the temporal evolution of fluid flow. Super-resolution can also work on spatial-temporal data. For example, STNet [47] addresses the spatial-temporal super-resolution of volumetric data.

*Method*. Spatial data inference typically include VAE, GAN and DRL methods.

- *Variational Autoencoder (VAE).* VAE [119] is a commonly used generative method that formulates generation as an autoregressive learning framework. The basic autoencoder architecture consists of an encoder extracting data features into latent representation vectors and a decoder reconstructing the data from the latent vectors. To allow GenAI to capture the variability in data, VAE builds on traditional autoencoder by modeling latent representation as probablistic instead of a fixed vector. During generation, the decoder sample from the distribution in the latent space and synthesize new data. VAE has been exploited for many data inference tasks in visualization such as data interpolation. For example, Latent Space Cartography [46] trains multiple VAE models with different hyperparameters on 24,000 emoji images. Users can explore the latent space of these VAEs and define customize semantic axes by selecting samples representing two ends of opposing concepts. Subsequently, linear interpolation is performed at constant intervals along the axis to generate intermediate samples showing the transitions of visual features of the emoji images.

- *Generative Adversarial Network (GAN).* GAN [120] models the process of generation as an adversarial learning framework, the basic components of which are the generator and the discriminator. The generator $G$ is designed to generate data that resembles real training data. The discriminator $D$ is designed to distinguish the data generated by the generator from the real data. The training alternates between the generator and discriminator to optimize a min-max problem with the following objective function. Different from the original GAN, spatial-temporal adversarial generation requires a spatial-temporal generator and discriminator. For example, STNet [47] builds a ConvLSTM structure for discriminator. Specifically, Convolution layers are used to extract spatial features. Then features of adjacent time steps are fed into ConvLSTM to evaluate temporal coherence. Global average pooling is used to produce the final single value score for realness.

- *Disentangled Representation Learning (DRL).* DRL with VAEs or GANs have been applied to visual analytics to identify interpretable dimensions interactively [48, 73, 69, 72, 74]. The disentangled dimensions can be subsequently controlled by users to generate meaningful data for augmentation. In the general literature of computer vision and
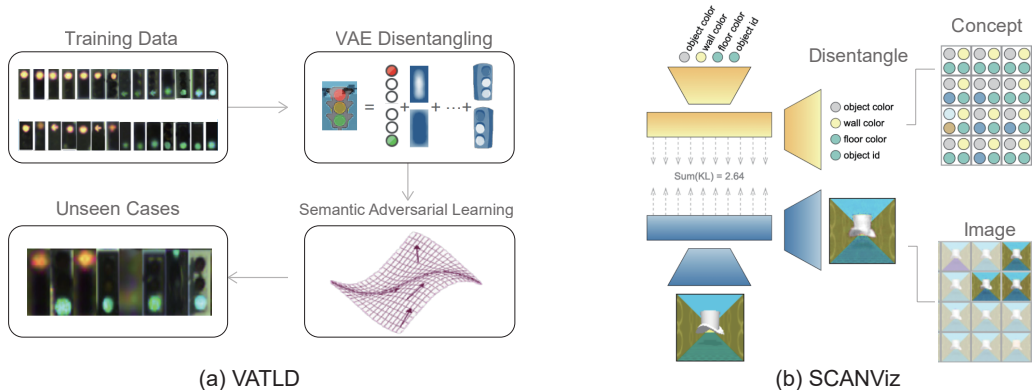
Figure 3: Data enhancement with disentangled representation learning, such as VATLD [48] and SCANViz [73].

computer graphics, disentanglement has been an essential technique for controllable generation [121, 122, 123]. A commonly used DRL architecture is $\beta$-VAE [124]. The objective function of $\beta$-VAE is a modification of the original VAE with an additional $\beta$ parameter. Experiments show that better chosen $\beta$ value (typically $< 1$) can produce more disentangled latent representation $\mathbf{z}$. For example, VATLD [48] is a visual analytics system that adapts $\beta$-VAE to extract user-interpretable features like colors, background and rotation from low-level features of traffic light images. With such interpretable features encoded in latent space, users can generate additional training examples in an interpretable manner to enhance the traffic light detection model. The DRL scheme distills potentially significant semantic dimensions in latent space representation for data summarization and semantic control by users. Particularly, two additional losses are introduced to the original $\beta$-VAE, namely the prediction loss and perceptual loss to ensure generation and reconstruction of more realistic traffic light images.

### 3.1.4. Discussion

GenAI methods like GANs are not designed to predict the accurate data values. Instead, they focus on generating reasonable data based on the given distribution of the real data. Such generation should not be overclaimed or misused for tasks that require accurate data features. A specific example is

outlier data, which might pose challenges as GenAI methods predominantly concentrate on learning the overall data distribution. Particularly, generative models aim to generate data that closely matches the majority of the training data. If the outliers are rare and significantly different from the majority of the data, the generative model may not capture them effectively. Outliers can be overlooked or underrepresented in the generated samples.

Despite being designed to embed data in more disentangled dimensions, the automatic DRL methods cannot guarantee the resulting dimensions are perfectly interpretable and disentangled. Consequently, the generative models built upon DRL are still largely a black box. For example, the choice of $\beta$ in $\beta$-VAE still remains largely heuristics. In addition, even though the dimensions extracted by DRL may be meaningful, it may not encode the visual properties users intend to explore, thus can potentially limit the customizable data exploration on users' part. Some recent work attempts to incorporate more user interaction to refine DRL through visual interface. For example, DRAVA [72] not only allows the meaning of DRL dimensions to be verified by users, but also enables user refinement. To facilitate user refinement of concept dimensions, they propose a light weight concept adaptor network on top of the VAE. The concept adaptor is a multi-class classifier to predict the correct grouping of data points along a selected dimensions for semantic clusterings. However, such interactions are still limited in some aspects. Users may only be able to verify and refine a small subset of dimensions, leaving many others unaddressed, because of the lack of overview for the relations between data points and different dimensions.

*3.2. Data Embedding*

Data embedding is an emergent technology that leverages GenAI to embed data into visualization images with information steganography. The data can be recovered losslessly from the visualization images, which are mostly 2D imagery data.

*3.2.1. Spatial Generation*

**Data**. Spatial generation in data embedding concerns QR code and visualization image data.

- *QR Code.* QR code data is a special type of data used as the visual coding scheme of the chart information such as metadata to be embedded, which can be processed together with the visualization images
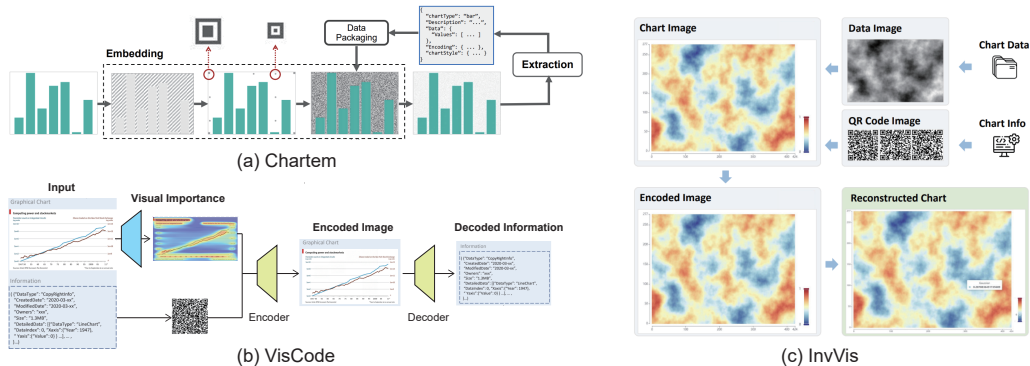
16

Figure 4: The data embedding pipeline of (a) Chartem [77], (b) VisCode [49] and (c) InvVis [50].

with neural networks. QR code is a reliable coding scheme allowing for error correction but avoiding artifacts in the encoded image [50, 49].

- *Visualization Image.* Although QR code can encode chart information, it has limited data encoding capacity. For the task of data embedding, visualization images are the essential medium to carry the encoded data and information. Large number of visualization images are needed to train the model. Synthetic data can be used for training, but real-world visualization image datasets such as VIS30K [125] and MASSVIS [126] have also been used to increase the generalizability and robustness of the model. To embed large quantities of underlying data for invertible visualization, data image that represents raw data produced by data-to-image (DTOI) method can also be used [50].

*Method*. The data embedding methods include boundary-aware segmentation network (BASNet) and invertible steganography network (ISN) models.

- *Saliency BASNet.* In order to assess the visual quality of the coded visualization image to ensure it is perceptually identical to the original image, it is insufficient to measure the pixel-wise mean square error because it neglects the varying importance of pixels across the visualization. To address this issue, VisCode [49] proposes a special visual importance network to predict the visual importance map for the chart image. Compared to traditional saliency-based method which

are applied to natural images but overlook the unique features of visualization images, the visual importance network can learn from eye-movement data of real users on chart images. Specifically, the model adopts a BASNet [127] architecture which is originally developed for salient object detection. The architecture is based on U-Net structure with residual blocks inspired by ResNet. The loss function combines BCEWIthLogits loss with structural similarity index (SSIM) to balance between segmentation accuracy and structural information.

- *Encoder-decoder ISN model.* To encode and decode the visualization image, QR code image and data image in one unified model, InvVis [50] introduces a concealing and revealing network. The concealing network and revealing network both consist of two major parts: feature fusion block (FFB) and invertible steganography network (ISN). FFB is designed to blend the features of data image and QR code image into visualization image while keeping minimal visual distortion. Specifically, FFB comprises four dense blocks and three common convolutional blocks. Dense blocks [128] are a special type of convolutional neural network architecture which contain multiple layer with dense connection (each layer is connected to all the preceding layers). Next, the ISN adds the invertible $1 \times 1$ convolution to the invertible neural network structure [129], which consist of several affine coupling layers. The authors also proposed using discrete wavelet transform (DWT) between the FFB and ISN to reduce texture-copying artifacts.

*3.2.2. Discussion*

Currently, the evaluation of the quality of data restoration is only limited to the data image, using generic metric such as root mean square error (RMSE). However, such pixel-wise metrics cannot fully reflect the accuracy of data restoration due to the absence of original data in the evaluation. In addition, the capacity of visualization image for data embedding is not infinite. Specifically, there can be an apparent trade-off between embedding capacity and image quality. To maintain image quality above certain threshold, it becomes more challenging to recover large amounts of data. To address this concern, more evaluation in practical scenario regarding precision requirements for the original data need to be conducted.

## 4. Visual Mapping Generation

For non-expert users, it is difficult to make appropriate visualization from data on their own for data analysis. GenAI plays an essential role in visual mapping synthesis for automatic visualization generation.

### 4.1. Visual Primitives Generation

The fundamental task of visual mapping generation is generating charts with the basic visual marks or visual primitives.

#### 4.1.1. Sequence Generation

*Data.* Sequence generation is often applied to visual mapping in different visualization grammars including concrete programming language and abstract code.

- *Visualization grammar.* Generative AI can be applied to generate visualization code in different grammars based on input data. For example, Vega-Lite code is a commonly used declarative visualization language [29, 130]. Data2Vis [29] treats generation of visual encoding as a sequence-to-sequence generation task which translate strings describing columns of data tables into Vega-Lite code sequences. Other online repositories such as Plotly also provides codes in other programming languages such as python. In addition, some studies generate abstract code instead of specific programming code. For example, Table2Charts [30], the authors define a more abstract chart template language including the essential visual elements and a set of grammar that summarizes the possible actions in the process of chart creation.

*Method.* Sequence generation used in visual mapping typically includes RNN, Deep Q Network, NL2VIS Transformer and the latest LLMs.

- *RNN-based Code Sequence Generation.* Some studies formulate the problem of visualization code generation as sequence-to-sequence generation. For example, Data2Vis [29] translates strings describing columns of data tables into Vega-Lite code sequences. For this task, the authors take inspiration from machine translation and adopt a encode-decoder architecture based on recurrent neural network model. Specifically, for the decoder, they construct a two-layer bidirectional RNN; for the decoder, another two-layer RNN is used to predict the next token in the

code sequence. Both the encoder and decoder leverage the Long Short Term Memory (LSTM) structure to enhance the model's ability to deal with longer sequence.

- *Deep Q network for encoding action prediction.* Some studies consider the task of visualization generation as the generation of abstract action tokens deciding key features of the charts, including data queries which select particular fields in data table and design choices which specify visual encoding operation. In this light, visualization generation can be formulated as an action prediction task which can be solved by Deep Q Network (DQN). For example, Table2Charts [30] develops a simple chart template language describing some essential actions for generation of six types of charts from data tables. According to this template, the authors construct a DQN for action prediction with a customized CopyNet architecture [131]. This network takes all the data fields and prefix action sequence as input and generate the next action token with a Gated Recurrent Unit (GRU) based RNN structure. In additon, to address the exposure bias problem with the previous teacher forcing training scheme which only learns the ground truth user generated results, Table2Charts adopts the search sampling approach of reinforcement learning to close the gap between training and inference.

- *Natural language to visualization models.* The Natural language to Visualization (NL2VIS) [78, 31, 79, 80] task can be formulated as follows: Given a natural language query ($NL$) over a dataset or relational database ($D$), the goal is to generate a visualization query (*e.g.* Vega-Lite) that is equivalent in meaning, valid for the specified $D$, and, when executed, will return a rendered visualization ($VIS$) result that aligns with the user's intent. For example, ADVISor [78] trained two separate neural networks to provide the NL2VIS functionality. Broadly, ADVISor's pipeline can be divided into two steps: (1) the NL2SQL step, and (2) the rule-based visualization generation step. ADVISor takes as input a NL question and data attributes associated with the datasets. Next, it first utilizes a BERT-based neural network to generate a vector representation ($q$) of the NL question and a corresponding header vector. Subsequently, the system utilizes an Aggregation network to deduce aggregate operators and a Data network to determine attribute selection and filtering conditions. Upon completing these steps, AD-

20

VISor first queries the dataset. It then maps the query results to a visualization based on a rule-based visualization algorithm The neural network modules within ADVISor are specifically trained to extract fragments of SQL queries from the given NL query, which indicates that ADVISor is not an end-to-end NL2VIS solution. On the contrary, ncNet [31] is an end-to-end NL2VIS method based on the Transformer [132] architecture. ncNet is trained on the first large-scale cross domain NL2VIS dataset nvBench [79]. ncNet utilizes a Transformer-based encoder-decoder framework, with both the encoder and decoder comprising self-attention blocks. This system accepts a NL query, a dataset, and an optional chart template as inputs. ncNet processes these inputs into embeddings and finally generates a flattened visualization query through an auto-regressive mechanism. Additionally, ncNet incorporates a visualization-aware decoding strategy, which allows for the generation of the final visualization query, with visualization-specific knowledge.

- *Large language model for visualization code generation.* Recently, some researchers realize the limitation of previous GenAI methods which only focus on a particular type of visualization code like Vega-Lite. To improve the flexibility of visualization code generation, some studies propose using large language model for more robust generation [32, 81, 83, 84, 85]. For example, LIDA [32] presents a pipeline called VIS-GENERATOR for AI generation of grammar-agnostic visualizations connecting data tables to generated visualizations with multiple steps. The VISGENERATOR consists of three sub-modules: code scaffold constructor, code generator and code executor. The code scaffold constructor generates code that imports language-specific dependencies like Matplotlib and constructs the empty function stub. Then, in code generator, taking as input the dataset summary and visualization goal, LLM is used in the fill-in-the-middle mode [133] to generate concrete visualization code of the given programming language. Finally, in the code executor, some filtering mechanisms such as self-consistency [134] and correctness probabilities [135] are incorporated to reduce errors. In the third step of LIDA, taking as input the dataset summary and visualization goal, LLM is used in the fill-in-the-middle mode [133] to generate concrete visualization code of different programming languages. Another recent study, LLM4Vis [81] proposes leveraging the in-context

learning ability of large language models to perform few-shot and zero-shot generation for the same design choice task as in VizML [39]. The key contribution of this method compared to previous supervised learning is that it reduces the need for large corpus of data-visualization pair training data and provides more explainable generation. The generation algorithm is retrieval-augmented with demonstration examples. First, data feature description is generated to enable GPT to take tabular dataset as input. Specifically, similar to VizML, with feature engineering as many as 120 single-column features and 80 cross-column features are extracted to represent the input dataset. These features are then serialized by TabLLM method [136], which utilizes a prompt to instruct ChatGPT to generate text description that elaborate on the feature values for each attribute. Next, the demonstration examples are selected from the training corpus by similarity retrieval based on feature description, which can fit into token limit of LLM without considering irrelevant samples. Subsequently, an iterative explanation generation bootstrapping module prompts LLM to not only predict the correct visual design choices but also generate explanation. Finally, all the relevant demonstration examples along with the explanation are fed into LLM to optimize in-context learning for generation of appropriate design choices for the input data. Recently some researchers also leverage LLM to refine the colormap [82].

### 4.1.2. Tabular Generation

*Data*. The visual mapping can also be simplified as predicting some tabular attributes such as design choices.

- *Design choices.* Instead of directly generating the visualization images, some studies focus on generating the most important design choices for the specifications of visualization, which is represented as tabular data with each attribute denoting one design dimension.

*Method*. The tabular generation methods of design choices include fully connected neural network for direct prediction and design parameter enumeration with AI scoring.

- *Fully connected neural network* can be combined with feature engineering of data, casting the generation problem into a prediction task. For

(a) Data2Vis          (b) DeepEye

Figure 5: Examples of visual primitive generation. (a) Data2Vis [29] adopts RNN-based sequence generation method. (b) DeepEye [17] combine design parameter enumeration with AI scoring.

example, VizML [39] builds a feed-forward neural network based on 841 dataset-level features extracted from input data table. The model predicts five design parameters of the appropriate visualizations with multi-head output layer, on both the encoding level and visualization level. For example, to generate the visualization type, one of the prediction head outputs a 6-class prediction scores for chart types including *Scatter, Line, Bar, Box, Histogram, Pie.*

- *Design parameters enumeration with AI scoring.* Some studies approach the generation problem from a different angle, using AI as the judge of the candidate generation results. For example, Deep-Eye [17, 86, 137, 138, 139] combines rule-based generation with data-driven machine learning to classify and rank meaningful visualizations. Specifically, based on a collected corpus of real world visualization cases, a classification model and a ranking model are trained. When users input a new dataset to visualize, the system first generates candidate visualizations by enumerating valid combinations of transformations and visual encoding in a pre-deifined search space. The classification model then determines whether a visualization candidate is meaningful. Subsequently, the ranking model sort the remaining meaningful visualization and recommend to users. In case the machine learning models do not yield satisfactory results, DeepEye also supports incorporation of expert-designed domain rules. The recommendations of

23

the machine learning model and the rule-based method can also be combined with a linear model. Similarly, Text-to-viz [18] adopt a hybrid method combining template-based enumeration with AI-powered comprehension of user inputs and relevance ranking.

### 4.1.3. Spatial Generation

*Data.* Spatial generation in visual mapping mainly concerns 2D sketch, density map and volumetric data.

- *Sketch.* Sketch is a special type of 2D image data with simple information of drawing trajectory. It has attracted much research interest in the broader field of generative AI [140] because it allows designers to follow their familiar workflow of prototype design and allows them to have spatial control over the generated results. Recently, some studies explore using sketch as a medium to facilitate fast prototyping of visualizations with sketch-to-vis generative AI [51, 141].

- *Density map.* Density maps are a special type of visualization that can vary dynamically depending on the time. Traditionally, the spatial temporal data collected for density map visualization is discrete and static. For smoother transition of density maps at different discrete observation times, some researchers propose to utilize generative AI [52].

- *Volumetric data.* Apart from the enhancement of volumetric data as introduced in Section 3.1.3, some studies also leverage GenAI methods for the rendering of such data [87, 88].

*Method.* Spatial generation methods for visual mapping mainly includes Faster R-CNN with validation model, GAN-based density map generation and GAN-based volume rendering.

- *Faster R-CNN for sketch recognition with validation model.* For the task of sketch to dashboard generation, latest AI-driven approach combines AI-powered chart recognition with rendering algorithms such as color palette recommendation and layout optimization. To detect the charts and the basic visual encoding features, LADV [51] first applies a Faster R-CNN network [142] which exploits region proposal to achieve
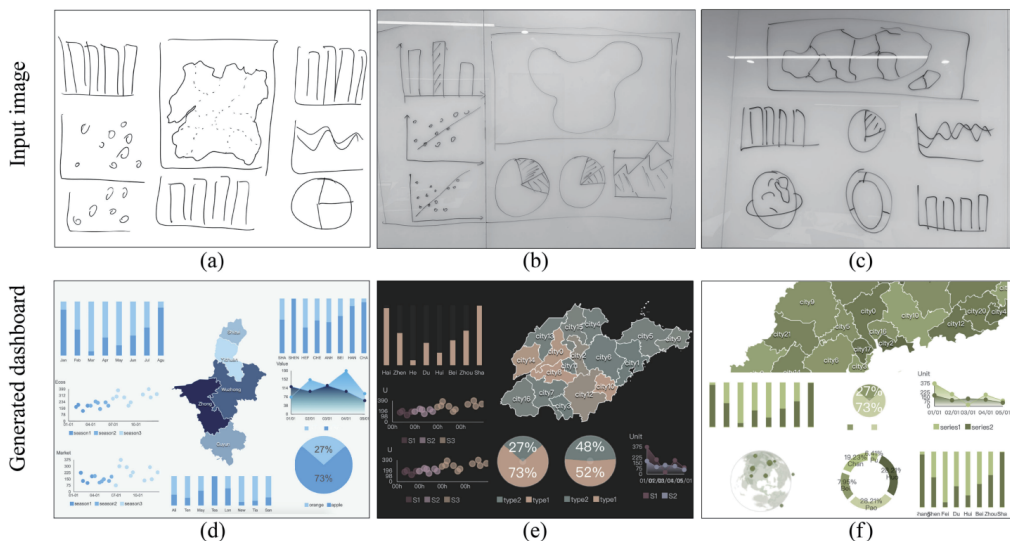
Figure 6: Examples of LADV [51] sketch to dashboard generation.

efficient object detection. Specifically, Faster R-CNN further accelerates Fast R-CNN by computing the region proposal with a deep CNN-based region proposal network which can share weights with the subsequent object detection network. In addition, to adapt Faster R-CNN to chart recognition, LADV [51] further incorporates a validation model to filter chart candidates, which trains a logistic regression for each chart type to capture the location and size.

- *GAN-based density generative model.* To generate dynamic density maps, density generative model is developed. For example, GenerativeMap [52] adopts a GAN-based generative model. Specifically, the authors first generate synthetic training data using Perlin noise. Then, they adapt Bidirectional Generative Adversarial Network (Bi-GAN) [143] to the density generation with enlarged convolution kernels and blocks inspired by ResNet [144] to enable processing of larger fields. Poisson blending is also used to make the density change more natural.

- *GAN-based Volume rendering.* Some researchers adopt GenAI for the rendering of 3D volumetric data [87, 88, 89, 90]. For example, Berger
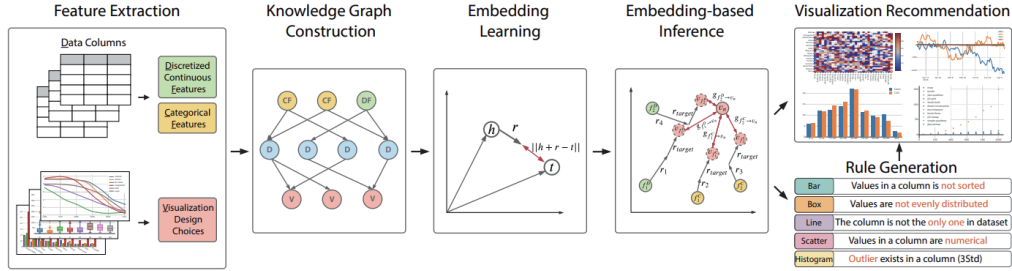
Figure 7: Knowledge-graph-empowered visual mapping generation in KG4VIS [42].

et al. [87] proposed a framework that combines two GANs for the task, namely opacity GAN and opacity-to-color translation GAN. Such approach breaks down the more difficult task of rendering volumes compared to the image generation task of the original GAN. Specifically, the opacity GAN learns to generate an opacity image given the input of viewpoint and opacity transfer function, which captures the shape, silhouette, and opacity. The second GAN translates the combined inputs of the viewpoint, the opacity transfer function's representation in the latent space, color transfer function, as well as the opacity image generated by the first GAN into the final rendered image.

### 4.1.4. Graph Generation

*Data*. The graph generation in visual mapping mainly involves knowledge graphs.

- *Knowledge graph.* Knowledge graph is a graph structure representation of knowledge that captures the relationships between entities in a particular domain, which is composed of nodes representing entities and edges representing relationships between these entities. Some studies also leverage knowledge graphs to support more explainable generation of visualization [42].

*Method*. The method for knowledge graph enhanced visual mapping generation is knowledge graph embedding.

- *Knowledge-graph embedding.* KG4VIS [42] applies knowledge graph to the design choice generation task of VizML [39]. The knowledge graph

26

is constructed with entities representing data features, data columns and visualization design choices, and the relations between them. Then, KG4Vis employs knowledge graph embedding method TransE [145] to learn neural embedding representations of all the entities and relations. When generating design choices for new input data, the model encode the input data in the embedding space and evaluate the scores of different candidate visual encoding rules with arithmetic operation adding relation vectors to entity vectors and measuring the distance from design choice vectors.

### 4.1.5. Discussion

**Limitations of the rules and the training data**. The training data can be a significant limiting factor for visual primitive generation. For example, Data2Vis [29] only trains the model on four major chart types, which significantly limits the scope and diversity of the generation. In addition, some hybrid methods combining rule-based components with GenAI may be limited by the rules that are not comprehensive enough. For example, the chart template defined in Table2Charts only [30] include the most basic visual mark and direct data field reference without even considering some basic operation like aggregation.

**Generation of visualization images vs. natural images**. Image generation has been a heavily researched topic in the general field of AI and computer vision. However, the generation of visualization images rarely adopts a fully end-to-end method without rule-based constraints. This is partly because of the difference between visualization images and natural images. Especially, compared to the irregular and complex visual features in natural images, the dominance of regular shapes rigid structures in visualization images makes it difficult for GenAI to accurately maintain. The reason is AI models are inherently stochastic and treats the image features as a whole with little knowledge of the structural constraints. To address this issue, future work may take inspiration from some recent computer vision studies that seek to incorporate structural information in image processing [146].

**3D diffusion**. Recently, inspired by the success of text-to-image generation diffusion model, some researchers seek to develop diffusion-based text-to-3D [16, 147] models to allow more intuitive interactive control of the generation, including natural language guided generation and editing. Such technology can potentially benefit generation of visualization, especially in providing multi-modal control.

**NL2VIS challenges**. Unquestionably, LLMs offer a complementary dimension to the NL2VIS system. However, the integration of LLMs into NL2VIS through prompt engineering presents certain limitations. First, relying solely on a simple prompting-based method may not effectively enable LLMs to fully comprehend the intricacies of the NL2VIS task. This approach could limit the model's ability to accurately interpret and respond to more complex visualization queries, potentially overlooking nuanced aspects of data visualization requirements. Second, current LLMs-based NL2VIS solutions often do not incorporate specific domain knowledge from the fields of visualization and data analysis into the LLMs. This absence of domain-specific integration can result in suboptimal performance, as the models may not leverage the rich contextual and technical knowledge necessary for producing highly accurate and relevant visualizations. Furthermore, current LLMs-based NL2VIS solutions struggle to guarantee the semantic correctness of generated visualizations, which is crucial for accurate data representation and interpretation. Additionally, these systems often face challenges in interactively fine-tuning results based on user feedback, a key aspect for achieving user-centered visualization design.

Given these challenges, future research directions may include: 1) Developing methods to integrate specific domain knowledge related to visualization and data analysis into LLMs. This could involve training models on specialized datasets or incorporating expert systems that guide the LLMs in understanding domain-specific visualization tasks, knowledge, and requirements. 2) Ensuring the semantic correctness of visualizations generated by LLMs, which can explore validation strategies that automatically check and confirm the accuracy and relevance of visualizations with respect to the underlying data. 3) Enhancing the interactivity of LLMs-based NL2VIS systems by incorporating more robust and flexible user feedback loops. For example, it can explore how to incorporate other modalities (*e.g.*clicks) user feedback. 4) Investigating hybrid models that combine the strengths of LLMs with traditional data visualization techniques. Such hybrid systems could leverage the natural language understanding capabilities of LLMs while ensuring adherence to best practices in data visualization.

## 5. Stylization

### 5.1. Style Transfer

Style broadly refers to the visual or aesthetic characteristics of an image, which oftentimes involves some global or overall features that can affect the viewers' general appreciation. More concretely, it can involve many aspects such as color, texture and layout. Some style transfer studies focus on the overall style while others focus on particular aspects like color.

Creating visualization in design practice relies on existing examples from the internet, offering inspirational visual materials toward style. This has spawned extensive research on transferring style attributes from these existing materials to facilitate the development of intended visualization designs. To transfer overall visual style, some research [148, 91, 54, 51] summarize style as a template and migrate these graphical attributes disentangled from the content to restyle new data source.

#### 5.1.1. Spatial Generation

*Data.* Researchers frequently employ visualizations in image format to train GenAI models, which are designed to extract specific attributes for transfer tasks.

- *Imagery data.* End-to-end GenAI models consume a substantial amount of visual images, primarily sourced from the internet or synthesized using tools like D3 [149]. Examples include MassVis [126], InfoVIF [94], and Visually29K [150].

- *Task-oriented Labeling.* General image datasets often lack paired attributes, necessitating an extraction stage for labeling task-oriented attributes. For instance, the color transfer task requires extracting the color map for subsequent training. Achieving the overall transfer necessitates considering multiple attributes to describe a comprehensive visual representation.

*Method.* The methods include color transfer and hybrid attribute transfer.

- *Color transfer.* As one of the most important visual channels in data visualization, generations of scholars have investigated the issue of color transfer [53, 151]. To extract the color at multiple scales in the Lab

histograms, Yuan et al. [53] employed a neural network featuring an atrous spatial pyramid structure, predicting the colormap of visualization and supporting discrete and continuous formats. Similarly, Huang et al. [91] approached the problem of color extraction in the foreground of visualization with Faster-RCNN. With the reference example as the natural image, some research [152, 92] generates a harmonious palette for a visualization based on color detection and extraction from images. For instance, Liu et al. [92] distinguished the salient subject in the image to extract the color with high visual importance that aligns with human perception.

- *Hybrid attribute transfer with Siamese Network.* Transferring multiple attributes from an example to the current design involves recognizing the content of the example and adapting it to the current design [54]. Lu et al. [94] curated a comprehensive infographic dataset and proposed a model based on YOLO to identify various visual elements, including text, icons, indices, and arrows. To maintain style consistency between the example and the current design, Vistylist [54] employs a Siamese Neural Network [153]. This network embeds visual elements into a 256-dimensional vector and compares the Euclidean distance between pairs. However, assessing the priority of different visual elements in a visualization poses a challenge. Huang et al. [91] proposed a restyling approach with an attention mechanism to weigh different visual properties for input visualizations. Accurate recognition of the example's content and the ability to reproduce it enable hybrid attribute transfer. This not only maintains style consistency but also generates a design tailored to the provided content. For example, when the template timeline has limited space, Chen et al. generated an extended timeline with a similar visual style to the template [93].

*5.1.2. Graph Generation*

*Data*. When considering the structure and imagery format of graphs, the data fed into GenAI models falls into two primary categories.

- *Graph Structure Data.* This includes nodes and edges in a graph, with the node feature vector and adjacency vector utilized to describe the graph structure that can be recognized by GenAI models. Various embedding techniques, such as node2vec [154], have been proposed to encode node information.

30

- *Imagery data.* Given that many graphs are in pixel format, GenAI models also process such data to extract low-level features for training.

`Method`. The method mainly includes graph layout transfer.

- *Graph layout transfer.* Some researchers use generative AI for graph layout transfer [43, 95, 44] which seeks to learn the style of graph layout from examples. For example, to help users intuitively produce diverse graph layouts from a given set of examples without mannually tweaking layout parameters, Kwon and Ma [43] designed a GenAI method based on encoder-decoder architecture combined with a 2D latent space. The model generally follows a VAE framework, taking as input graph layout features represented as relative pairwise distances and adjacency matrix. Then, GNN is used to compute graph-level representation of layout in both the encoder and decoder. In addition, the latent representation of layout $z_L$ is combined with node-level features through a fusion layer. The latent space is also visualized in a 2D map to facilitate exploration.

*5.1.3. Discussion*

**Domain knowledge**. Different visualizations have specific requirements and constraints for style transfer, often involving the integration of domain knowledge. For the categorical data, Zheng et al. [155] introduced a method to sample dominant colors from the image to preserve color discriminability, effectively enhancing and aiding in the interpretation of the patterns present. As for scientific visualization like terrain maps, domain-specific elements including continuity of elevation and hypsometric tints in aerial perspective are injected into the transfer process to convey the necessary information in a scientifically accurate manner [156].

**Reference image**. Furthermore, it is worth mentioning that the reference images used for style transfer are not limited to visualization examples alone. Significant works [157, 53, 91] have explored using natural images as reference sources for style transfer, taping into the inherent visual appeal and cognitive stimulation provided by natural images. Recent research has shown that natural images can also serve as an adorable source to stimulate human intelligence [92, 156].

## 5.2. Chart Embellishment

Visually embellished visualization showcases its memorability and expressiveness [158, 159, 160]. The creation of visually appealing and informative graphical enhancements necessitates design expertise. Fortunately, the advent of generative AI offers a strong framework to streamline the design process, particularly for pictorial visualization and glyph generation.

### 5.2.1. Spatial Generation

*Data.* The data that is concerned in GenAI for pictorial visualization is images with semantic correlation with the charts.

- *Image as pictorial embellishment.* In visualization, non-visualization images such as natural images or artistic images can be used as pictorial embellishments to enhance the visual appeal of the data being presented. Images can be added to charts, graphs, and other visualizations to provide additional context and meaning to the data [161, 162]. For example, an image of a product can be added to a sales chart to help viewers understand which product is being represented by each data point. Similarly, an image of a city skyline can be added to a map to help viewers identify the location being represented.

*Method.* The GenAI method for pictorial visualization mainly include Stable diffusion based techniques.

- *Stable diffusion.* Pictorial visualization plays a crucial role in seamlessly integrating semantic context into a chart. Instead of relying on pre-existing graphical elements sourced online and adjusting them to fit the desired visualization, generative AI takes an end-to-end approach by incorporating users' prompts. Recent advancements [32, 22, 96] in this field have led to the automation of the creation process through the transformation of chart components into semantic-related objects. For example, viz2viz [96] develops specific pipelines for generating various types of visualizations. The way they leverage textual prompts as input surpasses the previous language-oriented creation tools limited by a predetermined set of entities [18, 34], providing a robust semantic recognition to arbitrary user input. Furthermore, Xiao et al. [22] proposed a unified approach that classifies visual representations into

foreground and background. This approach provides users with an interface to select the intended data mask and incorporates an evaluation module to assess the visual distortion in the generated visualization.

### 5.2.2. Discussion

**Flexibility and controllability**. When comparing generative AI with the traditional methods of designing visual elements from scratch or retrieving relevant resources, generative AI offers several distinct advantages. It provides inspiration and eliminates the tedious and time-consuming process of adjusting elements to fit the data. Moreover, it empowers users by allowing them to personalize the style of the generated results with a simple utterance, saving significant time that would otherwise be expended on searching for appropriate resources across the vast expanse of the internet. Recent advancements in generative models, particularly in the field of text-to-image models, have achieved remarkable breakthroughs in enhancing control over the generated output, including layout [163], text content [164], vector graphics [165], etc. However, apart from general control, it is imperative to prioritize data integrity throughout the generation process, as visualizations serve to convey data patterns faithfully.

### 5.3. Textual Annotation

Textual annotation plays a pivotal role in the realm of data visualization, enhancing human interpretation, interaction, and comprehension. Text annotations incorporated in the visualization guide users' interactions with the artifact [166], explain what the data means [167], and prioritize certain interpretations of the data [168]. In this way, annotations act as cognitive aids, enhancing the overall user experience.

### 5.3.1. Sequence Generation

**Data.** The data in this task mainly includes text-integrated data and and contextually interpreted data.

- *Text-integrated data.* Text-integrated data refers to datasets that inherently come with accompanying textual content, such as articles, reports, or storytelling contexts. In these instances, the data is intertwined with text, forming a cohesive narrative or explanatory structure. Textual annotations in such datasets help make complex data more accessible and understandable to the audience, guide their attention to

trends that align with textual content, and provide a richer and more nuanced understanding of the data and its relevance to the text.

- *Contextually interpreted data.* For the datasets that lack accompanying textual narratives or descriptions, the challenge lies in analyzing the data to extract relevant insights and generate meaningful textual annotations that align with the visual elements. Effective textual annotations act as a bridge between complex datasets and audiences, providing a layer of interpretation that the raw data lacks on its own.

`Method.` To further automate the annotation generation process, researchers have developed innovative approaches, mainly through sequence generation.

- *Deep learning detection with template-based generation.* Lai et al. [169] employ a Mask R-CNN model to identify and extract visual elements in the target visualizations, along with their visual properties. The descriptive sentence is displayed beside the described focal areas as annotations. AutoCaption [35], a deep learning-powered scheme, generates captions for information charts by learning the noteworthy features aligned with human perception and leveraging one-dimensional residual neural networks to analyze relationships between visualization elements. These advances in automated annotation generation hold promise for applications in education and data overviews. Researchers in other areas such as NLP also study the problem of chart summarization [170].

- *Large Language Models (LLMs).* Recent developments in large language models (LLMs) have opened up new possibilities in generating engaging captions for generic data visualizations. Liew and Mueller [33] apply LLMs to produce descriptive captions for generic data visualizations like a scatterplot. Ko et al. [97] introduce a large language model (LLM) framework to generate rich and diverse natural language (NL) datasets using only Vega-Lite specifications as input. This underscores the growing role of prompt engineering techniques in shaping the future of annotation techniques and prompts a reevaluation of research directions.
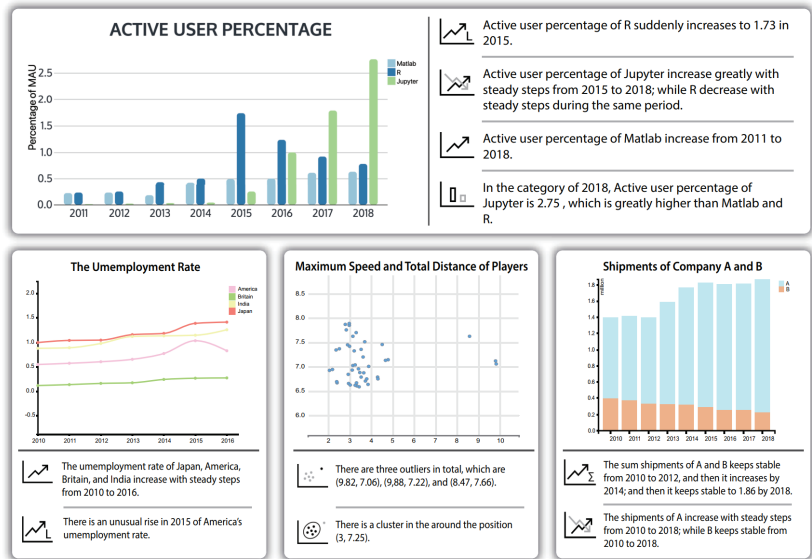
Figure 8: Annotated example cases of AutoCaption [35]. The chart type includes the bar chart, the scatterplot, and the line chart.

### 5.3.2. Discussion

**Challenges for textual annotation**. Researchers have consistently underscored the significance of annotations in visualization design, both at the visual memory level [126] and the cognitive level [171]. These studies reaffirm that annotations play an indispensable role in enhancing comprehension and retention of visual information [59]. In the realm of annotation for visualization, it is imperative to address three significant aspects for future research and development: mitigating occlusion problems, harnessing advanced techniques for automation, and enriching visual design. Firstly, a persistent challenge is the occlusion problem, wherein the annotations block the charts. Despite considerable efforts to improve the layout of annotated charts, this issue continues to hamper the effectiveness of annotations. Therefore, more suitable design space should be surveyed and innovative strategies for annotation placement should be considered. Secondly, recent advancements in deep learning, exemplified by large language models and diffusion models, offer remarkable potential for improving the efficiency of the automatic annotation process. These models can consider contextual information to produce annotations that are contextually relevant and strategically placed, alleviating the burden on users. Moreover, it is essential to enhance the design

35

of annotations by incorporating rich visual cues, which aid in highlighting patterns and facilitating deeper understanding. Researchers should ensure the visually engaging annotations complement the overall visual composition and improve user engagement.

### 5.4. Timeline & Story

Both timelines and storylines use lines to describe a sequence of events. Specifically, in a storyline visualization, each role is represented as a line.

### 5.4.1. Spatial Generation

*Data*. GenAI for story generation commonly concerns about optimizing the spatial layout of different story components.

- *Storyline Layout Data.* To training GenAI models to optimize the storyline layouts, it is necessary to generate a large number of high-quality storyline images. In PlotThread [55], Tang *et al.*construct a dataset of automatically-generated storyline layout pairs, consisting of a original layout and a optimized layout simulated by the optimization model with randomly-selected constraints.

*Method*. The method includes

- *Reinforcement Learning for image-based storyline.* To boost a collaborative design of storylines between AI and designers, Tang *et al.*[55] further proposed a reinforcement learning framework and introduced an authoring tool, PlotThread, that integrates an AI agent for efficient exploration and flexible customization. The goal of this is to imitate and improve users' intermediate results when optimizing storyling lines. Therefore, it is necessary to understand the states of different layouts, decompose a storyline into a sequence of interactive actions, and provide subsequent actions for layout optimization. They also define the reward as the similarity between the user layout and generated intermediate layout to improve the agent's prediction ability.

### 5.4.2. Graph Generation

*Data*. Researchers seek to understand visualization content automatically from the image input, which can be mainly categorized into raster image and vector image.

- *SVG.* A recent attempt is to apply GenAI technologies to vector charts due to the need of motioning images [100]. They use charts of SVG format to extract and model corresponding structural information between graphical elements without high computation cost.

*Method.*

- *Graph Nerual Network for structured dynamic charts.* Ling *et al.* [100] presented an automated method that transforms static charts into dynamic live charts for more effective communication and expressive presentation. To overcome the difficulty of generating dynamic live charts from static vector-based SVG, this study proposes using GNN for understanding chart and recovering data and visual encodings. Specifically, they first transform raw SVG into graph by a graph construction algorithm which extracts 5-dimensional node features including element type, node color, fill color, stroke color and stroke width; then it builds two types of edges including stroke-wise edges and element-wise edges. The constructed graph is then fed into two GNN-based encoder, each designed for one type of edges, to generate graph representations, which are subsequently passed to multi-layer perceptron to classify each graph element.

*5.4.3. Sequence Generation*

*Data.* To generate a complete story, most studies generate a sequence of data facts and ensemble them into a complete data story.

- *Relational data.* As the most basic format of visualization data, relational data is also a popular input of automatic story generation. GenAI are applied to generate textual descriptions for data tables [100] and construct the links between visuals and narrations through data table and word inputs [98].

- *Time Series Data.* Time series data is a widely used data type for a variety of visual analysis tasks, ranging from visual question answering to free exploration. The traditional tools are mostly designed for single-step guidance while GenAIs provide opportunities to build a continuous exploratory visual analysis process by extracting coherent data insights.

*Method.*

- *Large Language Model.* Ling *et al.* [100]'s work also leveraged large language models to create animated visuals and audio narrations, including narration with contextual information, narration with insights and narration rephrasing. To further enhance the interplay between visual animation and narration in data videos, Data Player [98] applies large language models to establish semantic connections between text and visualization and then recommends suitable animation presets with domain-knowledge constraints. Specifically, the authors design special prompt engineering with few-shot pre-defined examples illustrating how to output semantic links sequence provided the input of both data table and narration word.

- *Reinforcement learning.* Moreover, some researchers adopt RL-based sequence generation [34, 99, 101]. Shi *et al.* [99] build a reinforcement learning-based system to support the exploratory visual analysis of time series data. It constructs the agent's state and action space with domain knowledge to generate coherent data insights sequences as visual analysis recommendations. Specifically, the authors use the markovian decision process (MDP) model to formulate an EVA sequence as a sequence of state-action pairs. Then, the RL-based method seeks to maximize the cumulative reward, which combines familiarity reward and curiosity reward. In calculation of the curiosity reward, the casualCNN model is used to embedding time sequences of different length into equal lengths.

## 6. Interaction

GenAI methods such as large language models have demonstrated great potential for enhancing interaction in the broader field of human computer interaction [172, 9]. With GenAI methods, users can engage with the visualization charts, extracting novel insights and findings via a natural language interface, as known as Chart Question Answering. Given a collection of well-designed visualization charts, users can effortlessly navigate through this corpus to locate their desired chart using similarity search, which has also recently incorporated some GenAI techniques.

*6.1. Visualization Retrieval*

Having established a set of well-crafted visualization charts and share online, the subsequent question that arises is how we can help users in searching

for their desired visualizations within a given repository effectively and efficiently. This task is referred to as visualization retrieval. Engaging in visualization retrieval can offer significant advantages to several downstream tasks such as learning visualization design [173, 174], visualization reuse [45], visualization corpus construction [175, 56], web mining [176, 177, 178], and computational journalism [179].

Recently, some researchers adopt GenAI method to facilitate visualization retrieval tailored to user intent about visual structure or other features.

### 6.1.1. Spatial Generation

*Data*. Spatial generation methods for retrieval mainly involve raster visualization images.

- *Visualization images*. Recently, some studies leverage GenAI for enhanced representations of raster visualization images in retrieval, such as WYTIWYR [57] and LineNet [56].

*Method*. Methods for this task mainly include Triplet autoencoder and contrastive learning.

- *Triplet autoencoder*. LineNet [56] addresses the problem of line chart retrieval by considering both image-level and data-level similarity. For this purpose, a Triplet autoencoder is constructed with the backbone architecture of vision transformer [180]. Additionally, Luo et al. [56] also contribute a large-scale line chart corpus, named LineBench. This corpus contains over 115,000 line charts along with corresponding metadata from four real-world datasets, facilitating the study of similarity search in line chart visualizations.

- *Contrastive learning*. WYTIWYR [57] uses a contrastive language image pretraining (CLIP) [181] to facilitate zero-shot user intent alignment with visualization images.

### 6.1.2. Graph Generation

*Data*. The data mainly involves SVG format visualizations.

- *SVG*. Graph representation is also used by a recent study [45] to incorporate structural information in SVG-format visualization retrieval.

*Method*. The method mainly includes graph contrastive learning.

- *Graph contrastive learning.* Specifically, the InfoGraph [182] architecture is used, which is a contrastive graph representation learning model. Specifically, the model adopts GNN encoder and generate embedding vector for input graph structure, where the optimization scheme takes as input pairs of graphs and maximize the mutual information between one graph and its subgraph while minimizing mutual information with the subgraph in the other graph. Combining the graph representations with image-level visual representations, the visualization retrieval results prove to be more structurally consistent.

### 6.1.3. Discussion

**Retrieval augmented generation (RAG)**. In the field of GenAI, a recently popular topic is how to integrate retrieval into the generation pipeline to achieve knowledge-grounded generation and reduce the uncertainty of purely blackbox GenAI models [183, 184]. Recently some researchers also contemplate introducing such framework to visualization generation [80] to reduce task complexity and increase reliability of generated results. This work focuses on generating DV query sequence. However, the RAG framework has the potential to benefit many more different GenAI applications in visualization. For example, for infographics generation, we can take the best of both worlds by combining the previous retrieval-based methods [54, 161] with the latest purely GenAI methods [22]. In this way, users can benefit from both the reliable real-world examples and the creativity of GenAI models.

**Multi-modal composed retrieval**. WYTIWYR [57] introduces a retrieval prototype with the novel composed query which combines image input with text describing users' additional intent. Such multi-modal composed retrieval has attracted considerable attention in the general domain of image retrieval [185, 186, 187] and holds promise for improving the retrieval interaction for user intent alignment. Future work can further investigate the different combinations of multi-modal queries beyond text and image modalities and different logical compositions to allow for more flexible query interaction.

### 6.2. Chart Question Answering

In data analysis with information visualization, sometimes only providing the charts to the users is not adequate as it might be time-consuming for them

to comprehend complex information about the data in the chart. Chart question answering (CQA) [188, 189, 190, 108, 191, 192] is a burgeoning field of research which seeks to develop intelligent algorithms and systems to answer users' questions about the charts to expedite data analysis and enhance user interaction.

*6.2.1. Sequence Generation*

`Data`. CQA mainly considers chart questions data.

- *Chart questions.* Chart questions can be categorized according to different attributes [192], including factual/open-ended, visual/non-visual and simple/compositional. In addition, in a more general sense, there can be different modalities of input as query about the chart.

`Method`. The CQA methods mainly include chart elements detection and vision-language model.

- *Chart elements detection.* Many AI-powered CQA models rely on detection of chart elements and structure to facilitate extraction of relevant information from visualization as explicit prior for generation of answers, such as PlotQA [102], FigureNet [103], DVQA [104], STL-CQA [36] and LEAF-QA [105]. For example, PlotQA [102] utilizes both Visual Element Detection (VED) and Object Character Recognition (OCR) to extract key information from charts. In visual element detection, Faster R-CNN is used while in the OCR a traditional method is adopted. Subsequently, the extracted chart elements is converted into a knowledge graph, which is combined with log-linear ranking of logical forms extracted from the question with compositional semantic parsing to generate the answer. Other works adopt more complex model fully based on neural network. For example, DVQA [104] develops a multi-output model that is capable of answering both generic questions and chart specific questions. The model contains an OCR sub-network composed of a CNN-based bounding box predictor and a GRU-based character-level decoder to extract the text. Based on the results of OCR sub-network, DVQA improves the Stacked Attention Network (SAN) [193] for general visual question answering with additional dynamic encoding, which can adapt to chart specific vocabulary. Although some most recent works still [106] rely on Chart elements
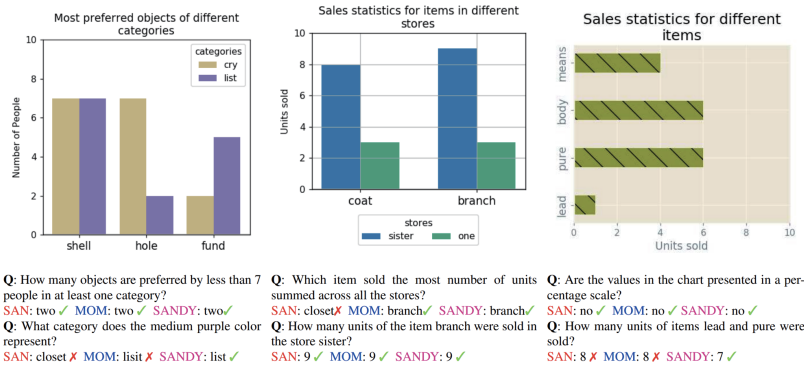
41

Figure 9: Chart question answering examples in DVQA [104].

detection and data extraction like ChartOCR [194] as an integral part, a few study [108] start to utilize a slightly different type of algorithm, which is OCR-free document image understanding such as Donut [195]. Donut adopts the Swin Transformer architecture and is pretrained on an OCR-pseudo task. However, in the inference stage it does not require external explicit OCR information and simply acts like an image encoder.

- *Vision-language model for multi-modal fusion.* With the growing power of generative AI, especially the multi-modal feature fusion capability, some recent studies simplifies the chart question answering pipeline with unified vision-language model, such as ChartQA [106], PReFIL [107], Unichart [108] and ChartLlama [37]. For example, ChartQA [106] builds a baseline model utilizing VL-T5 [196], a pre-trained unified vision-language model for text generation conditioned on multi-modal inputs. ChartQA also proposes their own model, VisionTaPas, which is a multi-modal extension of the TaPas [197] model. The original TaPas model is designed for answering questions over table, where a table is flattened into a sequence of words, converting the problem into essentially a unimodal-input text generation task. For this task, the BERT [198] architecture is extended with additional embeddings to represent table structure and context, including embeddings of segment, row/column, rank, and previous answer. In the VisionTaPas model, a Vision Transformer (ViT) [199] model is utilized to extract chart image features into embeddings, as ViT has proven to be more powerful

42

than CNN in many vision tasks. Next, a cross-modality encoder is constructed to fuse the multi-modal embeddings of ViT and TaPas, combining information of both text and chart images for end-to-end generation of answers. Some recent research also explores instruction tuning of pre-trained large vision-language model for more flexible generation of answer [200].

*6.2.2. Discussion*

**Other modalities**. Most existing CQA systems only consider single modality natural language input as the primary means of interaction with visualization. However, other studies have shown the importance of interactions in other modalities, such as body movement [201], touch and pen [202] and gesture [203]. A few studies have attempted to achieve multi-modal inputs by combining natural language or speech with mouse, pen or touch interactions [204, 205, 206]. Nevertheless, these tentative works rely on traditional rule-based methods for quick prototyping and have not exploited the latest GenAI methods for CQA as introduced above. To achieve multi-modal CQA, data-driven generative AI promises more flexibility than traditional methods in diverse real-world scenarios. For example, some researchers have started to utilized GPT to support sketch interaction for generation of chart findings documentation [207].

**Combination with data embedding**. As shown in our introduction of the algorithms, most GenAI-based CQA methods still depend on explicit detection of chart elements and underlying data for generation of precise answers. Such detection may not be always accurate and robust for complex real-world visualization images due to the diverse styles and visualization types as well as additional noises. One possible strategy to circumvent this issue is combining the data embedding method introduced in Section 3.2 with CQA. With additional information about the data embedded in the chart images, the performance of CQA can be expected to further improve.

**The promise of more precise vision-language model**. Recent development of vision-language model is showing a trend of higher precision for more fine-grained detail in the images. Segment Anything [208] Model can locate specific semantic segment in the image given users visual or textual prompt. Similarly, Grounding-Dino [209] can even more accurately generate bounding box for particular objects in images with users' prompts. In addition, LlaVa [210] allows users to flexibly ask about different levels of image contents from overall features to details. For the task of CQA which requires

so much precision that additional detection models are needed, these powerful vision-language models have the potential to significantly simplify the pipeline, leading towards a universal multi-modal model for chart interaction.

**GenAI for visual analytics**. Going beyond the interaction with single charts, some researchers recently are exploring the possibility of extending GenAI, particularly large language models to more complicated visual analytics workflow [109, 211]. For example, LEVA [109] utilizes LLM to assist multiple stages of visual analytics including onboarding, exploration, and summarization. With the development of LLM agent technology [212], GenAI can potentially take the role of humans in some visual analytics tasks. It is an important question for future research to define new paradigm for human-AI collaboration in data visualization and analysis.

## 7. Research Challenges And Opportunities

### 7.1. Evaluating GenAI for Visualization

The increasing use of GenAI in the production of complex and creative visualizations. Given the essential role of rigorous evaluation in visualization design, it becomes crucial to apply similar assessment standards to AI-generated visualizations. The distinct characteristics and challenges presented by AI-driven visualization processes necessitate careful adaptation of evaluation metrics and methodologies. While traditional metrics such as efficiency [21] and aesthetic [159] remain fundamental in evaluating AI-generated visualizations, the advent of AI techniques introduce additional, specific metrics that must be considered. From the migration of assessment metrics for GenAI, the following assessment metrics are likely to be considered for evaluating differnet application of GenAI in visualization.

- **Accuracy and Fidelity.** Ensuring accuracy and fidelity in AI-generated visualizations is paramount, particularly when applying stylization techniques. Techniques such as semantic contextualization in visualization [22] face the challenge of balancing data integrity with aesthetic appeal. This is crucial because real-life objects often do not conform to the rigid outlines typical in model-generated images, posing a risk to the accuracy of the visual representation.

- **Intent Alignment and Controllbility.** This criterion assesses the degree to which AI-generated visualizations align with the user's intent

44

Table 3: Examples of GenAI4VIS datasets.

| Dataset | Data Format | Source | Supported Tasks |
|---|---|---|---|
| VizNet [217] | Real world tables | Web-crawled | Data inference |
| VIS30K [125] | Chart images | Extracted from papers | Data embedding |
| Data2Vis [29] | Table-code pairs | Synthetic | Table2VIS generation |
| nvBench [79] | NL-code pairs | Synthetic | NL2VIS generation |
| MV [218] | MV-layout labels | Extracted from papers | Layout transfer |
| Chart-to-text [170] | Chart-text pairs | Web-crawled | Text annotation |
| Beagle [219] | SVG-type labels | Web-crawled | Visualization retrieval |
| LineBench [56] | Chart-data pairs | Synthesis with annotation | Visualization retrieval |
| PlotQA [102] | Chart-QA pairs | Crowd-sourcing + Synthetic | CQA |

and their ability to influence the outcome. In Natural Language Interaction (NLI), controllability pertains to the user's efficacy in steering the output of language-based AI systems during iterative interactions [213, 214]. Furthermore, in end-to-end generation processes, it is essential that the AI-generated visualizations are sensitive to and align with the user's specific requirements, such as style or query objectives.

- **Robustness and Consistency.** Lastly, evaluating the robustness and consistency of AI-generated visualizations across different scenarios is a key metric, ensuring reliability and applicability in diverse contexts [215]. Regarding LLM, it may blending fact with fiction and generating non-factual content, which called hallucination problem [216]. For example, in doing vqa tasks, especially in domains with specific requirements for accuracy, evaluating the hallucinations of the generated content is essential.

- **Bias and Ethics.** The potential biases inherent in AI algorithms and the ethical implications of their outputs necessitate careful examination.The generative model may face potential criticism on copyright or bias issues, as the training process digests a huge amount of data obtained from the web, which is unfiltered and imbalanced.

In short, the field needs to update evaluation methods and criteria continually to keep pace with advancing GenAI technologies in assessing AI-generated visualizations.

*7.2. Dataset*

As GenAI is data-driven, these methods heavily depends on the training data. Indeed, most previous works applying GenAI to visualization build

their own dataset or utilize the datasets created by prior works [23]. Even in the era of large language models which are pre-trained on much larger general purpose dataset, a domain-specific visualization dataset can serve as valuable reference and knowledge base for efficient prompting and improving the reliability of GenAI results. The quality, quantity, and diversity of the dataset thus have a significant impact on the generative performance and the output quality, as it determines how the GenAI model perceives and understands the patterns and semantics of the generation requirement and generated content.

In this regard, several aspects warrant special attention in future research. First, the diversity is important, as a diverse training dataset helps the AI model learn a broader range of topics, styles and other design patterns in real-world visualization. This diversity enables the model to generate content that is more versatile and contextually appropriate in different situations. However, many datasets used in training GenAI4VIS models are less diverse than real-world data [220], partly because most researchers collect or synthesize their training data for prototyping of their generative methods, without sufficient consideration of more complex authentic cases. Therefore, building on existing GenAI4VIS studies, one important direction for future improvement is understanding the lack of diversity in current training datasets and supplement them accordingly.

Second, the heterogeneous data in different formats can reduce the reusability. For instance, visualizations data are in a wide range of forms including raster images, SVG and different types of codes like Vega-Lite and Python. This discrepancy necessitates the curation of datasets in different formats anew for different generation tasks. This can also be a significant limiting factor for the size of the dataset because similar data in other formats cannot be utilized. This in turn may lead to overfitting and other difficulties in training GenAI for visualization. To address this issue, more robust visualization retargeting methods need to be developed to align different formats of data, such as translating Vega-Lite code to Python code and extracting graphical SVG structure from raster visualization images. For example, recently some researchers have been exploring the idea of using large language models to generate various annotations for visualization datasets [97].

In addition, we provide examples of some existing datasets that can be applied to different GenAI4VIS tasks, as shown in Table 3. We can find that some researchers seek to address the lack of GenAI4VIS dataset by either collecting and transforming real-world data or synthesizing data. For example,

due to the lack of NL2VIS benchmarks, nvBench [79, 221] proposes utilizing the existing NL2SQL benchmark dataset which can be transformed into NL2VIS benchmark. Furthermore, we can see that some larger scale real-world datasets such as VIS30K [125] and VizNet [217] so far can only facilitate data enhancement tasks in GenAI4VIS because of the lack of annotations about the visual mapping process. As we mentioned above, how to strike a balance between synthetic methods with high scalability and real-world data collection which requires more manual effort or a complex retargeting process can be a critical issue. LineBench [56] is a large-scale line chart visualization corpus (with 115,000 line charts) with the associated source dataset, the underlying data D for rendering, and the rendered visualization V in the form of an image, which can facilitate the study of similarity search of line chart visualizations. In this light, the perspective of data-centric explainable AI [222, 223] is particularly relevant. Many visual analytics studies for explainable AI seek to help users explore the models from the data perspective to gain insights about the potential biases, yet most of these works look at general-purpose AI or GenAI models. In other words, there is not enough self-reflection studies from the visualization domain to diagnose GenAI4VIS models along with their training data when most researchers are rushing to apply GenAI to various subtasks in the visualization pipeline.

### 7.3. GenAI4VIS vs. Generative Visualization

In the area of digital art, there is some distinction between AI art and generative art (or algorithmic art), where the latter term largely refers to traditional procedural generation algorithms with rule-based or optimization-based methods such as graph grammar or genetic algorithm [224]. In contrast, AI art mostly encompasses end-to-end purely deep-learning-based generation methods such as GAN, VAE or diffusion models. However, in the field of visualization, there is little discussion about such distinction. In many GenAI4VIS studies, researchers often introduce a hybrid approach, integrating many rule-based constraints and procedures with partially AI-powered methods. For example, VizML [39] incorporates more than 800 hand-crafted features in the input, while restricting the generation output to predicting only a few basic visual structures. In essence, visualization has been relying on grammar-based generation which explicitly prescribe the mapping from data to visual structures and views with a suite of different codes and rules such as Vega-Lite and visual design guidelines. To some extent, this can limit the effort to fully harness the power of GenAI, mainly because the models

47

cannot directly learn the distribution of the final rendered images conditioned on data and user input, which is ultimately what visualization presents to users. This is vastly different from more mature GenAI applications in other areas. For example, in spatial generation, latest GenAI technology can skip most traditional image synthesis and 3D modeling procedures and directly render the 2D images or 3D models. LIDA [32] makes an early effort towards a more integrated GenAI4VIS pipeline. However, LIDA's pipeline is still divided into separate sequence generation for visualization code and spatial generation for visualization stylization in a linear workflow, where the two GenAI models do not share knowledge. One problem due to this disconnection in LIDA, for example, is that the stylization stage cannot maintain the accuracy of the visual structure with respect to the data because the image-based stable diffusion model in the second stage is completely ignorant of structural information in the previous sequence generation. Moreover, some recent studies in AI show that merging two large pretrained models using techniques like knowledge distillation can not only produce a versatile merged model but also boost the performance for downstream tasks that require knowledge from both models [225], which provides inspiration for potential strategies to improve integration of GenAI4VIS models.

In fact, the gap between visualization and GenAI pipelines is not necessarily a downside, as this signifies opportunities for future research to combine the advantages while mitigating the respective disadvantages. On the one hand, visualization researchers can think about how to directly model the mapping between data and views in the end-to-end statistical learning framework of GenAI, which can provide more effective learning and evaluation based on the final visual representations. For example, this means that raster visualization images from real world sources can also be directly utilized as training data as long as it is annotated with user requirement text labels, without needing further complex explicit chart element extraction for retargeting. In this way, visualization can potentially harness the true power of multi-modal GenAI based on large pretrained vision-language models, which is showing more accurate control down to the pixel-level in recent research [226]. On the other hand, the intermediate operations like those in visualization should not be discarded by the GenAI4VIS pipeline because GenAI can be more explainable and controllable if users are allowed to inspect and intervene in the key intermediate steps. However, such intervention should not be the superficial rule-based constraints in hybrid methods. Instead, researchers can take inspiration from GenAI research such as Con-

trolNet [227] and LayoutDiffusion [228] to embed the control into the model itself. Alternatively, researchers can develop interactive tools to support human intervention in the generation process [229].

## 8. Conclusion

The burgeoning GenAI technology is promising for applications in the visualization domain. Because of GenAI's impressive capacity to model the transformation and design process by learning from real data, it can benefit a range of visualization tasks like data enhancement, visual mapping generation, stylization and interaction. Different types of GenAI methods have been applied to these tasks due to different data structures, including sequence generation, tabular generation, spatial generation and graph generation. With the advent of latest GenAI technology like large language model and diffusion model, new opportunities emerge to revolutionize GenAI4VIS methods. However, task-specific challenges still exist due to the unique characteristics of visualization tasks, which demands further investigation. Moreover, general challenges in evaluation and datasets require some rethinking about the GenAI4VIS pipeline beyond simply borrowing state-of-the art GenAI methods. We hope this survey can help researchers reflect on existing GenAI4VIS research from a technical perspective and provide some inspiration for future research opportunities, with a vision for improved integration of GenAI in visualization.

## References

[1] A. Key, B. Howe, D. Perry, C. Aragon, VizDeck: self-organizing dashboards for visual analytics, in: Proceedings of the ACM SIGMOD International Conference on Management of Data, 2012, pp. 681–684.

[2] S. Zhu, G. Sun, Q. Jiang, M. Zha, R. Liang, A survey on automatic infographics and visualization recommendations, Visual Informatics 4 (3) (2020) 24–40.

[3] A. Wu, Y. Wang, X. Shu, D. Moritz, W. Cui, H. Zhang, D. Zhang, H. Qu, AI4VIS: Survey on artificial intelligence approaches for data visualization, IEEE Transactions on Visualization and Computer Graphics 28 (12) (2021) 5049–5070.

[4] Q. Wang, Z. Chen, Y. Wang, H. Qu, A survey on ML4VIS: Applying machine learning advances to data visualization, IEEE Transactions on Visualization and Computer Graphics 28 (12) (2021) 5134–5153.

[5] C. Wang, J. Han, Dl4SciVis: A state-of-the-art survey on deep learning for scientific visualization, IEEE Transactions on Visualization and Computer Graphics 29 (8) (2023) 3714–3733.

[6] X. Qin, Y. Luo, N. Tang, G. Li, Making data visualization more efficient and effective: a survey, VLDB J. 29 (1) (2020) 93–117.

[7] M. Abukmeil, S. Ferrari, A. Genovese, V. Piuri, F. Scotti, A survey of unsupervised generative models for exploratory data analysis and representation learning, ACM Computing Surveys 54 (5) (2021) 1–40.

[8] J. Gui, Z. Sun, Y. Wen, D. Tao, J. Ye, A review on generative adversarial networks: Algorithms, theory, and applications, IEEE Transactions on Knowledge and Data Engineering 35 (4) (2023) 3313–3332.

[9] Y. Hou, M. Yang, H. Cui, L. Wang, J. Xu, W. Zeng, C2Ideas: Supporting creative interior color design ideation with large language model, arXiv preprint arXiv:2401.12586 (2024).

[10] S. Xiao, L. Wang, X. Ma, W. Zeng, TypeDance: Creating semantic typographic logos from image through personalized generation, arXiv preprint arXiv:2401.11094 (2024).

[11] R. Huang, H.-C. Lin, C. Chen, K. Zhang, W. Zeng, PlantoGraphy: Incorporating iterative design process into generative artificial intelligence for landscape rendering, arXiv preprint arXiv:2401.17120 (2024).

[12] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, B. Ommer, High-resolution image synthesis with latent diffusion models, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2022, pp. 10684–10695.

[13] A. Ramesh, P. Dhariwal, A. Nichol, C. Chu, M. Chen, Hierarchical text-conditional image generation with clip latents, arXiv preprint arXiv:2204.06125 (2022).

[14] OpenAI, GPT-4 technical report (2023). `arXiv:2303.08774`.

[15] H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal, E. Hambro, F. Azhar, et al., Llama: Open and efficient foundation language models, arXiv preprint arXiv:2302.13971 (2023).

[16] B. Poole, A. Jain, J. T. Barron, B. Mildenhall, DreamFusion: Text-to-3d using 2d diffusion, in: The International Conference on Learning Representations, 2022.

[17] Y. Luo, X. Qin, N. Tang, G. Li, DeepEye: Towards automatic data visualization, in: Proceedings of the IEEE International Conference on Data Engineering, 2018, pp. 101–112.

[18] W. Cui, X. Zhang, Y. Wang, H. Huang, B. Chen, L. Fang, H. Zhang, J.-G. Lou, D. Zhang, Text-to-viz: Automatic generation of infographics from proportion-related natural language statements, IEEE Transactions on Visualization and Computer Graphics 26 (1) (2019) 906–916.

[19] J. Mackinlay, P. Hanrahan, C. Stolte, Show me: Automatic presentation for visual analysis, IEEE Transactions on Visualization and Computer Graphics 13 (6) (2007) 1137–1144.

[20] M. Chen, D. Ebert, H. Hagen, R. S. Laramee, R. Van Liere, K.-L. Ma, W. Ribarsky, G. Scheuermann, D. Silver, Data, information, and knowledge in visualization, IEEE Computer Graphics and Applications 29 (1) (2008) 12–19.

[21] E. R. Tufte, The visual display of quantitative information, Vol. 2, Graphics press Cheshire, CT, 2001.

[22] S. Xiao, S. Huang, Y. Lin, Y. Ye, W. Zeng, Let the chart spark: Embedding semantic context into chart with text-to-image generative model, IEEE Transactions on Visualization and Computer Graphics 30 (1) (2024) 284–294.

[23] C. Chen, Z. Liu, The state of the art in creating visualization corpora for automated chart analysis, Computer Graphics Forum 42 (3) (2023) 449–470.

[24] C. R. Qi, H. Su, K. Mo, L. J. Guibas, PointNet: Deep learning on point sets for 3d classification and segmentation, in: Proceedings of the IEEE

conference on computer vision and pattern recognition, 2017, pp. 652–660.

[25] S. K. Card, J. Mackinlay, B. Shneiderman, Readings in information visualization: using vision to think, Morgan Kaufmann, 1999.

[26] L. McNabb, R. S. Laramee, Survey of surveys (SoS)-mapping the landscape of survey papers in information visualization 36 (3) (2017) 589–617.

[27] L. Shen, E. Shen, Y. Luo, X. Yang, X. Hu, X. Zhang, Z. Tai, J. Wang, Towards natural language interfaces for data visualization: A survey, IEEE Transactions on Visualization and Computer Graphics 29 (6) (2023) 3121–3144.

[28] M. Tennekes, E. de Jonge, Tree colors: color schemes for tree-structured data, IEEE Transactions on Visualization and Computer Graphics 20 (12) (2014) 2072–2081.

[29] V. Dibia, Ç. Demiralp, Data2vis: Automatic generation of data visualizations using sequence-to-sequence recurrent neural networks, IEEE Computer Graphics and Applications 39 (5) (2019) 33–46.

[30] M. Zhou, Q. Li, X. He, Y. Li, Y. Liu, W. Ji, S. Han, Y. Chen, D. Jiang, D. Zhang, Table2Charts: Recommending charts by learning shared table representations, in: Proceedings of the ACM SIGKDD Conference on Knowledge Discovery & Data Mining, 2021, pp. 2389–2399.

[31] Y. Luo, N. Tang, G. Li, J. Tang, C. Chai, X. Qin, Natural language to visualization by neural machine translation, IEEE Transactions on Visualization and Computer Graphics (2021) 1–1.

[32] V. Dibia, LIDA: A tool for automatic generation of grammar-agnostic visualizations and infographics using large language models, in: Proceedings of the Annual Meeting of the Association for Computational Linguistics, 2023, pp. 113–126.

[33] A. Liew, K. Mueller, Using large language models to generate engaging captions for data visualizations, arXiv preprint arXiv:2212.14047 (2022).

[34] D. Shi, X. Xu, F. Sun, Y. Shi, N. Cao, Calliope: Automatic visual data story generation from a spreadsheet, IEEE Transactions on Visualization and Computer Graphics 27 (2) (2020) 453–463.

[35] C. Liu, L. Xie, Y. Han, D. Wei, X. Yuan, AutoCaption: An approach to generate natural language description from visualization automatically, in: IEEE Pacific Visualization Symposium (PacificVis), 2020, pp. 191–195.

[36] H. Singh, S. Shekhar, STL-CQA: Structure-based transformers with localization and encoding for chart question answering, in: Proceedings of the Conference on Empirical Methods in Natural Language Processing, 2020, pp. 3275–3284.

[37] Y. Han, C. Zhang, X. Chen, X. Yang, Z. Wang, G. Yu, B. Fu, H. Zhang, ChartLlama: A multimodal llm for chart understanding and generation, arXiv preprint arXiv:2311.16483 (2023).

[38] N. Park, M. Mohammadi, K. Gorde, S. Jajodia, H. Park, Y. Kim, Data synthesis based on generative adversarial networks, Proceedings of the VLDB Endowment 11 (10) (2018) 1071–1083.

[39] K. Hu, M. A. Bakker, S. Li, T. Kraska, C. Hidalgo, Vizml: A machine learning approach to visualization recommendation, in: Proceedings of the CHI Conference on Human Factors in Computing Systems, 2019, pp. 1–12.

[40] N. Shi, J. Xu, S. W. Wurster, H. Guo, J. Woodring, L. P. Van Roekel, H.-W. Shen, GNN-Surrogate: A hierarchical and adaptive graph neural network for parameter space exploration of unstructured-mesh ocean simulations, IEEE Transactions on Visualization and Computer Graphics 28 (6) (2022) 2301–2313.

[41] Y. Zhang, J. Li, C. Xu, Graph-based latent space traversal for new molecules discovery, in: Proceedings of the International Symposium on Visual Information Communication and Interaction, 2023, pp. 1–8.

[42] H. Li, Y. Wang, S. Zhang, Y. Song, H. Qu, KG4Vis: A knowledge graph-based approach for visualization recommendation, IEEE Transactions on Visualization and Computer Graphics 28 (1) (2021) 195–205.

[43] O.-H. Kwon, K.-L. Ma, A deep generative model for graph layout, IEEE Transactions on Visualization and Computer Graphics 26 (1) (2019) 665–675.

[44] S. Song, C. Li, Y. Sun, C. Wang, VividGraph: Learning to extract and redesign network graphs from visualization images, IEEE Transactions on Visualization and Computer Graphics 29 (7) (2023) 3169–3181.

[45] H. Li, Y. Wang, A. Wu, H. Wei, H. Qu, Structure-aware visualization retrieval, in: Proceedings of the CHI Conference on Human Factors in Computing Systems, 2022, pp. 1–14.

[46] Y. Liu, E. Jun, Q. Li, J. Heer, Latent space cartography: Visual analysis of vector space embeddings, Computer Graphics Forum 38 (3) (2019) 67–78.

[47] J. Han, H. Zheng, D. Z. Chen, C. Wang, STNet: An end-to-end generative framework for synthesizing spatiotemporal super-resolution volumes, IEEE Transactions on Visualization and Computer Graphics 28 (1) (2021) 270–280.

[48] L. Gou, L. Zou, N. Li, M. Hofmann, A. K. Shekar, A. Wendt, L. Ren, VATLD: A visual analytics system to assess, understand and improve traffic light detection, IEEE Transactions on Visualization and Computer Graphics 27 (2) (2020) 261–271.

[49] P. Zhang, C. Li, C. Wang, VisCode: Embedding information in visualization images using encoder-decoder network, IEEE Transactions on Visualization and Computer Graphics 27 (2) (2020) 326–336.

[50] Y. L. Huayuan Ye, Chenhui Li, C. Wang, InvVis: Large-scale data embedding for invertible visualization, IEEE Transactions on Visualization and Computer Graphics 30 (1) (2024) 1139–1149.

[51] R. Ma, H. Mei, H. Guan, W. Huang, F. Zhang, C. Xin, W. Dai, X. Wen, W. Chen, LADV: Deep learning assisted authoring of dashboard visualizations from images and sketches, IEEE Transactions on Visualization and Computer Graphics 27 (9) (2020) 3717–3732.

[52] C. Chen, C. Wang, X. Bai, P. Zhang, C. Li, Generativemap: Visualization and exploration of dynamic density maps via generative learning

model, IEEE Transactions on Visualization and Computer Graphics 26 (1) (2019) 216–226.

[53] L.-P. Yuan, W. Zeng, S. Fu, Z. Zeng, H. Li, C.-W. Fu, H. Qu, Deep colormap extraction from visualizations, IEEE Transactions on Visualization and Computer Graphics 28 (12) (2021) 4048–4060.

[54] Y. Shi, P. Liu, S. Chen, M. Sun, N. Cao, Supporting expressive and faithful pictorial visualization design with visual style transfer, IEEE Transactions on Visualization and Computer Graphics 29 (1) (2022) 236–246.

[55] T. Tang, R. Li, X. Wu, S. Liu, J. Knittel, S. Koch, T. Ertl, L. Yu, P. Ren, Y. Wu, Plotthread: Creating expressive storyline visualizations using reinforcement learning, IEEE Transactions on Visualization and Computer Graphics 27 (2) (2020) 294–303.

[56] Y. Luo, Y. Zhou, N. Tang, G. Li, C. Chai, L. Shen, Learned data-aware image representations of line charts for similarity search, Proceedings of the ACM on Management of Data 1 (1) (2023) 88:1–88:29.

[57] S. Xiao, Y. Hou, C. Jin, W. Zeng, WYTIWYR: A user intent-aware framework with multi-modal inputs for visualization retrieval, Computer Graphics Forum 42 (3) (2023) 311–322.

[58] J. Xia, J. Li, S. Chen, H. Qin, S. Liu, A survey on interdisciplinary research of visualization and artificial intelligence, Scientia Sinica (Informationis) 51 (2021) 1777–1801.

[59] Q. Chen, S. Cao, J. Wang, N. Cao, How does automation shape the process of narrative visualization: A survey of tools, IEEE Transactions on Visualization and Computer Graphics (2023) 1–20.

[60] Y. He, S. Cao, Y. Shi, Q. Chen, K. Xu, N. Cao, Leveraging large models for crafting narrative visualization: A survey, arXiv preprint arXiv:2401.14010 (2024).

[61] S. Di Bartolomeo, V. Schetinger, J. L. Adams, A. M. McNutt, M. El-Assady, M. Miller, Doom or deliciousness: Challenges and opportunities for visualization in the age of generative models, Computer Graphics Forum 42 (3) (2023) 423–435.

[62] W. Yang, M. Liu, Z. Wang, S. Liu, Foundation models meet visualizations: Challenges and opportunities, Computational Visual Media (2023).

[63] J. Fan, T. Liu, G. Li, J. Chen, Y. Shen, X. Du, Relational data synthesis using generative adversarial networks: A design space exploration, Proceedings of the VLDB Endowment 13 (11) (2020) 1962–1975.

[64] H. Chen, S. Jajodia, J. Liu, N. Park, V. Sokolov, V. Subrahmanian, FakeTables: Using gans to generate functional dependency preserving tables with bounded real data., in: Proceedings of IJCAI, 2019, pp. 2074–2080.

[65] X. Qinl, C. Chai, N. Tang, J. Li, Y. Luo, G. Li, Y. Zhu, Synthesizing privacy preserving entity resolution datasets, in: Proceedings of IEEE International Conference on Data Engineering, 2022, pp. 2359–2371.

[66] L. Xu, M. Skoularidou, A. Cuesta-Infante, K. Veeramachaneni, Modeling tabular data using conditional gan, in: Proceedings of the International Conference on Neural Information Processing Systems, Vol. 32, 2019.

[67] J. Zhang, G. Cormode, C. M. Procopiuc, D. Srivastava, X. Xiao, PrivBayes: Private data release via bayesian networks, ACM Transactions on Database Systems (TODS) 42 (4) (2017) 1–41.

[68] Y. Wang, Z. Zhong, J. Hua, DeepOrganNet: on-the-fly reconstruction and visualization of 3d/4d lung models from single-view projections by deep deformation network, IEEE Transactions on Visualization and Computer Graphics 26 (1) (2019) 960–970.

[69] W. He, L. Zou, A. K. Shekar, L. Gou, L. Ren, Where can we help? a visual analytics approach to diagnosing and improving semantic segmentation of movable objects, IEEE Transactions on Visualization and Computer Graphics 28 (1) (2021) 1040–1050.

[70] Z. Zhou, Y. Hou, Q. Wang, G. Chen, J. Lu, Y. Tao, H. Lin, Volume upscaling with convolutional neural networks, in: Proceedings of the Computer Graphics International Conference, 2017, pp. 1–6.

[71] S. Wiewel, M. Becher, N. Thuerey, Latent space physics: Towards learning the temporal evolution of fluid flow, Computer Graphics Forum 38 (2) (2019) 71–82.

[72] Q. Wang, S. L'Yi, N. Gehlenborg, DRAVA: Aligning human concepts with machine learning latent dimensions for the visual exploration of small multiples, in: Proceedings of the CHI Conference on Human Factors in Computing Systems, 2023, pp. 1–15.

[73] J. Wang, W. Zhang, H. Yang, SCANViz: Interpreting the symbol-concept association captured by deep neural networks through visual analytics, in: IEEE Pacific Visualization Symposium (PacificVis), 2020, pp. 51–60.

[74] N. Evirgen, X. Chen, GANravel: User-driven direction disentanglement in generative adversarial networks, in: Proceedings of the CHI Conference on Human Factors in Computing Systems, 2023, pp. 1–15.

[75] H. Singh, N. McCarthy, Q. U. Ain, J. Hayes, ChemoVerse: Manifold traversal of latent spaces for novel molecule discovery, arXiv preprint arXiv:2009.13946 (2020).

[76] W. Zheng, J. Li, Y. Zhang, Desirable molecule discovery via generative latent space exploration, Visual Informatics 7 (4) (2023) 13–21.

[77] J. Fu, B. Zhu, W. Cui, S. Ge, Y. Wang, H. Zhang, H. Huang, Y. Tang, D. Zhang, X. Ma, Chartem: reviving chart images with data embedding, IEEE Transactions on Visualization and Computer Graphics 27 (2) (2020) 337–346.

[78] C. Liu, Y. Han, R. Jiang, X. Yuan, Advisor: Automatic visualization answer for natural-language question on tabular data, in: IEEE Pacific Visualization Symposium (PacificVis), IEEE, 2021, pp. 11–20.

[79] Y. Luo, N. Tang, G. Li, C. Chai, W. Li, X. Qin, Synthesizing Natural Language to Visualization (NL2VIS) Benchmarks from NL2SQL Benchmarks, in: Proceedings of the International Conference on Management of Data, p. 1235–1247.

[80] Y. Song, X. Zhao, R. C.-W. Wong, D. Jiang, RGVisNet: A hybrid retrieval-generation neural framework towards automatic data visualization generation, in: Proceedings of the ACM SIGKDD Conference on Knowledge Discovery and Data Mining, 2022, pp. 1646–1655.

[81] L. Wang, S. Zhang, Y. Wang, E.-P. Lim, Y. Wang, LLM4Vis: Explainable visualization recommendation using ChatGPT, arXiv preprint arXiv:2310.07652 (2023).

[82] C. Shi, W. Cui, C. Liu, C. Zheng, H. Zhang, Q. Luo, X. Ma, NL2Color: Refining color palettes for charts with natural language, IEEE Transactions on Visualization and Computer Graphics 30 (1) (2024) 814–824.

[83] S. Li, X. Chen, Y. Song, Y. Song, C. Zhang, Prompt4Vis: Prompting large language models with example mining and schema filtering for tabular data visualization, arXiv preprint arXiv:2402.07909 (2024).

[84] Y. Tian, W. Cui, D. Deng, X. Yi, Y. Yang, H. Zhang, Y. Wu, Chart-GPT: Leveraging llms to generate charts from abstract natural language, IEEE Transactions on Visualization and Computer Graphics (2024) 1–15.

[85] G. Li, X. Wang, G. Aodeng, S. Zheng, Y. Zhang, C. Ou, S. Wang, C. H. Liu, Visualization generation with large language models: An evaluation, arXiv preprint arXiv:2401.11255 (2024).

[86] Y. Luo, X. Qin, N. Tang, G. Li, X. Wang, Deepeye: Creating good data visualizations by keyword search, in: Proceedings of the International Conference on Management of Data, 2018, pp. 1733–1736.

[87] M. Berger, J. Li, J. A. Levine, A generative model for volume rendering, IEEE Transactions on Visualization and Computer Graphics 25 (4) (2018) 1636–1650.

[88] F. Hong, C. Liu, X. Yuan, DNN-VolVis: Interactive volume visualization supported by deep neural network, in: IEEE Pacific Visualization Symposium (PacificVis), 2019, pp. 282–291.

[89] L. Wu, J. Y. Lee, A. Bhattad, Y.-X. Wang, D. Forsyth, DIVeR: Real-time and accurate neural radiance fields with deterministic integration

for volume rendering, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2022, pp. 16200–16209.

[90] W. Gan, H. Xu, Y. Huang, S. Chen, N. Yokoya, V4d: Voxel for 4d novel view synthesis, IEEE Transactions on Visualization and Computer Graphics (2023) 1–14.

[91] D. Huang, J. Wang, G. Wang, C.-Y. Lin, Visual style extraction from chart images for chart restyling, in: Proceedings of the International Conference on Pattern Recognition, 2021, pp. 7625–7632.

[92] S. Liu, M. Tao, Y. Huang, C. Wang, C. Li, Image-driven harmonious color palette generation for diverse information visualization, IEEE Transactions on Visualization and Computer Graphics (2022) 1–16.

[93] Z. Chen, Y. Wang, Q. Wang, Y. Wang, H. Qu, Towards automated infographic design: Deep learning-based auto-extraction of extensible timeline, IEEE Transactions on Visualization and Computer Graphics 26 (1) (2019) 917–926.

[94] M. Lu, C. Wang, J. Lanir, N. Zhao, H. Pfister, D. Cohen-Or, H. Huang, Exploring visual information flows in infographics, in: Proceedings of the CHI conference on human factors in computing systems, 2020, pp. 1–12.

[95] Y. Wang, Z. Jin, Q. Wang, W. Cui, T. Ma, H. Qu, DeepDrawing: A deep learning approach to graph drawing, IEEE Transactions on Visualization and Computer Graphics 26 (1) (2019) 676–686.

[96] J. Wu, J. J. Y. Chung, E. Adar, viz2viz: Prompt-driven stylized visualization generation using a diffusion model, arXiv preprint arXiv:2304.01919 (2023).

[97] H.-K. Ko, H. Jeon, G. Park, D. H. Kim, N. W. Kim, J. Kim, J. Seo, Natural language dataset generation framework for visualizations powered by large language models, arXiv preprint arXiv:2309.10245 (2023).

[98] L. Shen, Y. Zhang, H. Zhang, Y. Wang, Data player: Automatic generation of data videos with narration-animation interplay, IEEE Transactions on Visualization and Computer Graphics 30 (1) (2024) 109–119.

[99] Y. Shi, B. Chen, Y. Chen, Z. Jin, K. Xu, X. Jiao, T. Gao, N. Cao, Supporting Guided Exploratory Visual Analysis on Time Series Data with Reinforcement Learning, IEEE Transactions on Visualization and Computer Graphics (2023) 1–11.

[100] L. Ying, Y. Wang, H. Li, S. Dou, H. Zhang, X. Jiang, H. Qu, Y. Wu, Reviving static charts into live charts, arXiv preprint arXiv:2309.02967 (2023).

[101] G. Wu, S. Guo, J. Hoffswell, G. Y.-Y. Chan, R. A. Rossi, E. Koh, Socrates: Data story generation via adaptive machine-guided elicitation of user feedback, IEEE Transactions on Visualization and Computer Graphics 30 (1) (2024) 131–141.

[102] N. Methani, P. Ganguly, M. M. Khapra, P. Kumar, PlotQA: Reasoning over scientific plots, in: Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision, 2020, pp. 1527–1536.

[103] R. Reddy, R. Ramesh, A. Deshpande, M. M. Khapra, FigureNet: A deep learning model for question-answering on scientific plots, in: International Joint Conference on Neural Networks, IEEE, 2019, pp. 1–8.

[104] K. Kafle, B. Price, S. Cohen, C. Kanan, DVQA: Understanding data visualizations via question answering, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2018, pp. 5648–5656.

[105] R. Chaudhry, S. Shekhar, U. Gupta, P. Maneriker, P. Bansal, A. Joshi, Leaf-QA: Locate, encode & attend for figure question answering, in: Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision, 2020, pp. 3512–3521.

[106] A. Masry, X. L. Do, J. Q. Tan, S. Joty, E. Hoque, ChartQA: A benchmark for question answering about charts with visual and logical reasoning, in: Findings of the Association for Computational Linguistics, 2022, pp. 2263–2279.

[107] K. Kafle, R. Shrestha, S. Cohen, B. Price, C. Kanan, Answering questions about data visualizations using efficient bimodal fusion, in: Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision, 2020, pp. 1498–1507.

[108] A. Masry, P. Kavehzadeh, X. L. Do, E. Hoque, S. Joty, UniChart: A universal vision-language pretrained model for chart comprehension and reasoning, in: Proceedings of the Conference on Empirical Methods in Natural Language Processing, 2023, pp. 14662–14684.

[109] Y. Zhao, Y. Zhang, Y. Zhang, X. Zhao, J. Wang, Z. Shao, C. Turkay, S. Chen, LEVA: Using large language models to enhance visual analytics, IEEE Transactions on Visualization and Computer Graphics (2024) 1–17.

[110] R. Gómez-Bombarelli, J. N. Wei, D. Duvenaud, J. M. Hernández-Lobato, B. Sánchez-Lengeling, D. Sheberla, J. Aguilera-Iparraguirre, T. D. Hirzel, R. P. Adams, A. Aspuru-Guzik, Automatic chemical design using a data-driven continuous representation of molecules, ACS Central Science 4 (2) (2018) 268–276.

[111] Z. Wu, S. Pan, F. Chen, G. Long, C. Zhang, S. Y. Philip, A comprehensive survey on graph neural networks, IEEE Transactions on Neural Networks and Learning Systems 32 (1) (2020) 4–24.

[112] T. White, Sampling generative networks, arXiv preprint arXiv:1609.04468 (2016).

[113] J. Y. Yen, Finding the k shortest loopless paths in a network, Management Science 17 (11) (1971) 712–716.

[114] A. Radford, L. Metz, S. Chintala, Unsupervised representation learning with deep convolutional generative adversarial networks, arXiv preprint arXiv:1511.06434 (2015).

[115] Y. Li, B. Sixou, F. Peyrin, A review of the deep learning methods for medical images super resolution problems, Innovation and Research in BioMedical Engineering 42 (2) (2021) 120–133.

[116] S.-W. Huang, C.-T. Lin, S.-P. Chen, Y.-Y. Wu, P.-H. Hsu, S.-H. Lai, AugGan: Cross domain adaptation with gan-based data augmentation, in: Proceedings of the European Conference on Computer Vision, 2018, pp. 718–731.

[117] J. Choi, T. Kim, C. Kim, Self-ensembling with gan-based data augmentation for domain adaptation in semantic segmentation, in: Proceedings of the IEEE/CVF International Conference on Computer Vision, 2019, pp. 6830–6840.

[118] X. Liu, Y. Zou, L. Kong, Z. Diao, J. Yan, J. Wang, S. Li, P. Jia, J. You, Data augmentation via latent space interpolation for image classification, in: Proceedings of the International Conference on Pattern Recognition (ICPR), 2018, pp. 728–733.

[119] D. P. Kingma, M. Welling, Auto-encoding variational bayes, arXiv preprint arXiv:1312.6114 (2013).

[120] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, Y. Bengio, Generative adversarial networks, Communications of the ACM 63 (11) (2020) 139–144.

[121] L. Tran, X. Yin, X. Liu, Disentangled representation learning gan for pose-invariant face recognition, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2017, pp. 1415–1424.

[122] I. Jeon, W. Lee, M. Pyeon, G. Kim, Ib-GAN: Disentangled representation learning with information bottleneck generative adversarial networks, in: Proceedings of the AAAI Conference on Artificial Intelligence, Vol. 35, 2021, pp. 7926–7934.

[123] C. Zhou, F. Zhong, C. Öztireli, CLIP-PAE: Projection-augmentation embedding to extract relevant features for a disentangled, interpretable and controllable text-guided face manipulation, in: Proceedings of the ACM SIGGRAPH Conference, 2023, pp. 1–9.

[124] C. P. Burgess, I. Higgins, A. Pal, L. Matthey, N. Watters, G. Desjardins, A. Lerchner, Understanding disentangling in beta-vae, arXiv preprint arXiv:1804.03599 (2018).

[125] J. Chen, M. Ling, R. Li, P. Isenberg, T. Isenberg, M. Sedlmair, T. Möller, R. S. Laramee, H.-W. Shen, K. Wünsche, et al., Vis30k: A collection of figures and tables from ieee visualization conference publications, IEEE Transactions on Visualization and Computer Graphics 27 (9) (2021) 3826–3833.

[126] M. A. Borkin, A. A. Vo, Z. Bylinskii, P. Isola, S. Sunkavalli, A. Oliva, H. Pfister, What makes a visualization memorable?, IEEE Transactions on Visualization and Computer Graphics 19 (12) (2013) 2306–2315.

[127] X. Qin, Z. Zhang, C. Huang, C. Gao, M. Dehghan, M. Jagersand, BAS-Net: Boundary-aware salient object detection, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2019, pp. 7479–7489.

[128] G. Huang, Z. Liu, L. Van Der Maaten, K. Q. Weinberger, Densely connected convolutional networks, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2017, pp. 4700–4708.

[129] L. Dinh, D. Krueger, Y. Bengio, Nice: Non-linear independent components estimation, arXiv preprint arXiv:1410.8516 (2014).

[130] A. Narechania, A. Srinivasan, J. Stasko, NL4DV: A toolkit for generating analytic specifications for data visualization from natural language queries, IEEE Transactions on Visualization and Computer Graphics 27 (2) (2020) 369–379.

[131] J. Gu, Z. Lu, H. Li, V. O. Li, Incorporating copying mechanism in sequence-to-sequence learning, in: Proceedings of the Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), 2016, pp. 1631–1640.

[132] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, I. Polosukhin, Attention is all you need, in: Proceedings of the International Conference on Neural Information Processing Systems, 2017.

[133] M. Bavarian, H. Jun, N. Tezak, J. Schulman, C. McLeavey, J. Tworek, M. Chen, Efficient training of language models to fill in the middle, arXiv preprint arXiv:2207.14255 (2022).

[134] X. Wang, J. Wei, D. Schuurmans, Q. Le, E. Chi, S. Narang, A. Chowdhery, D. Zhou, Self-consistency improves chain of thought reasoning in language models, arXiv preprint arXiv:2203.11171 (2022).

[135] S. Kadavath, T. Conerly, A. Askell, T. Henighan, D. Drain, E. Perez, N. Schiefer, Z. Hatfield-Dodds, N. DasSarma, E. Tran-Johnson, et al.,

Language models (mostly) know what they know, arXiv preprint arXiv:2207.05221 (2022).

[136] S. Hegselmann, A. Buendia, H. Lang, M. Agrawal, X. Jiang, D. Sontag, TabLLM: Few-shot classification of tabular data with large language models, in: International Conference on Artificial Intelligence and Statistics, PMLR, 2023, pp. 5549–5581.

[137] Y. Luo, X. Qin, C. Chai, N. Tang, G. Li, W. Li, Steerable self-driving data visualization, IEEE Transactions on Knowledge and Data Engineering 34 (1) (2020) 475–490.

[138] X. Qin, Y. Luo, N. Tang, G. Li, DeepEye: An automatic big data visualization framework, Big Data Mining and Analytics 1 (1) (2018) 75–82.

[139] X. Qin, Y. Luo, N. Tang, G. Li, Deepeye: Visualizing your data by keyword search, in: EDBT, OpenProceedings.org, 2018, pp. 441–444.

[140] A. Voynov, K. Aberman, D. Cohen-Or, Sketch-guided text-to-image diffusion models, in: ACM SIGGRAPH Conference Proceedings, 2023, pp. 1–11.

[141] Z. Teng, Q. Fu, J. White, D. C. Schmidt, Sketch2Vis: Generating data visualizations from hand-drawn sketches with deep learning, in: Proceedings of the IEEE International Conference on Machine Learning and Applications (ICMLA), 2021, pp. 853–858.

[142] S. Ren, K. He, R. Girshick, J. Sun, Faster R-CNN: Towards real-time object detection with region proposal networks, in: Proceedings of the International Conference on Neural Information Processing Systems, Vol. 28, 2015.

[143] J. Donahue, P. Krähenbühl, T. Darrell, Adversarial feature learning, arXiv preprint arXiv:1605.09782 (2016).

[144] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recognition, in: Proceedings of the Conference on Computer Vision and Pattern Recognition, 2016, pp. 770–778.

[145] Z. Sun, Z.-H. Deng, J.-Y. Nie, J. Tang, RotatE: Knowledge graph embedding by relational rotation in complex space, in: International Conference on Learning Representations, 2018.

[146] K. Han, Y. Wang, J. Guo, Y. Tang, E. Wu, Vision GNN: An image is worth graph of nodes, in: Proceedings of the International Conference on Neural Information Processing Systems, Vol. 35, 2022, pp. 8291–8303.

[147] E. Sella, G. Fiebelman, P. Hedman, H. Averbuch-Elor, Vox-E: Text-guided voxel editing of 3d objects, in: Proceedings of the IEEE/CVF International Conference on Computer Vision, 2023, pp. 430–440.

[148] J. Harper, M. Agrawala, Converting basic d3 charts into reusable style templates, IEEE Transactions on Visualization and Computer Graphics 24 (3) (2017) 1274–1286.

[149] M. Bostock, V. Ogievetsky, J. Heer, $D^3$ data-driven documents, IEEE Transactions on Visualization and Computer Graphics 17 (12) (2011) 2301–2309.

[150] Z. Bylinskii, N. W. Kim, P. O'Donovan, S. Alsheikh, S. Madan, H. Pfister, F. Durand, B. Russell, A. Hertzmann, Learning visual importance for graphic designs and data visualizations, in: Proceedings of the Annual ACM symposium on user interface software and technology, 2017, pp. 57–69.

[151] L.-P. Yuan, Z. Zhou, J. Zhao, Y. Guo, F. Du, H. Qu, InfoColorizer: Interactive recommendation of color palettes for infographics, IEEE Transactions on Visualization and Computer Graphics 28 (12) (2021) 4252–4266.

[152] Y. Shi, S. Chen, P. Liu, J. Long, N. Cao, Colorcook: Augmenting color design for dashboarding with domain-associated palettes, Proceedings of the ACM on Human-Computer Interaction 6 (CSCW2) (2022) 1–25.

[153] M. Lagunas, E. Garces, D. Gutierrez, Learning icons appearance similarity, Multimedia Tools and Applications 78 (2019) 10733–10751.

[154] A. Grover, J. Leskovec, node2vec: Scalable feature learning for networks, in: Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 2016, pp. 855–864.

[155] Q. Zheng, M. Lu, S. Wu, R. Hu, J. Lanir, H. Huang, Image-guided color mapping for categorical data visualization, Computational Visual Media 8 (4) (2022) 613–629.

[156] M. Wu, Y. Sun, S. Jiang, Adaptive color transfer from images to terrain visualizations, IEEE Transactions on Visualization and Computer Graphics (2023) 1–16.

[157] J. Poco, A. Mayhua, J. Heer, Extracting and retargeting color mappings from bitmap images of visualizations, IEEE Transactions on Visualization and Computer Graphics 24 (1) (2017) 637–646.

[158] R. Borgo, A. Abdul-Rahman, F. Mohamed, P. W. Grant, I. Reppa, L. Floridi, M. Chen, An empirical study on using visual embellishments in visualization, IEEE Transactions on Visualization and Computer Graphics 18 (12) (2012) 2759–2768.

[159] L. Harrison, K. Reinecke, R. Chang, Infographic aesthetics: Designing for the first impression, in: Proceedings of the ACM Conference on Human Factors in Computing Systems, 2015, pp. 1187–1190.

[160] S. Haroz, R. Kosara, S. L. Franconeri, Isotype visualization: Working memory, performance, and engagement with pictographs, in: Proceedings of ACM Conference on Human Factors in Computing Systems, 2015, pp. 1191–1200.

[161] D. Coelho, K. Mueller, Infomages: Embedding data into thematic images, Computer Graphics Forum 39 (3) (2020) 593–606.

[162] J. E. Zhang, N. Sultanum, A. Bezerianos, F. Chevalier, Dataquilt: Extracting visual elements from images to craft pictorial visualizations, in: Proceedings of the CHI Conference on Human Factors in Computing Systems, 2020, pp. 1–13.

[163] J. Cheng, X. Liang, X. Shi, T. He, T. Xiao, M. Li, Layoutdiffuse: Adapting foundational diffusion models for layout-to-image generation, arXiv preprint arXiv:2302.08908 (2023).

[164] J. Chen, Y. Huang, T. Lv, L. Cui, Q. Chen, F. Wei, Textdiffuser: Diffusion models as text painters, arXiv preprint arXiv:2305.10855 (2023).

[165] P. Zhang, N. Zhao, J. Liao, Text-guided vector graphics customization, arXiv preprint arXiv:2309.12302 (2023).

[166] E. Segel, J. Heer, Narrative visualization: Telling stories with data, IEEE Transactions on Visualization and Computer Graphics 16 (6) (2010) 1139–1148.

[167] A. Cairo, The Functional Art: An introduction to information graphics and visualization, New Riders, 2012.

[168] J. Hullman, N. Diakopoulos, Visualization rhetoric: Framing effects in narrative visualization, IEEE Transactions on Visualization and Computer Graphics 17 (12) (2011) 2231–2240.

[169] C. Lai, Z. Lin, R. Jiang, Y. Han, C. Liu, X. Yuan, Automatic annotation synchronizing with textual description for visualization, in: Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems, 2020, pp. 1–13.

[170] S. Kantharaj, R. T. Leong, X. Lin, A. Masry, M. Thakkar, E. Hoque, S. Joty, Chart-to-Text: A large-scale benchmark for chart summarization, in: Proceedings of the Annual Meeting of the Association for Computational Linguistics, 2022, pp. 4005–4023.

[171] S. Latif, Z. Zhou, Y. Kim, F. Beck, N. W. Kim, Kori: Interactive synthesis of text and charts in data documents, IEEE Transactions on Visualization and Computer Graphics 28 (1) (2021) 184–194.

[172] Z. Wang, L. Yuan, L. Wang, B. Jiang, Z. Wei, VirtuWander: Enhancing multi-modal interaction for virtual tour guidance through large language models, arXiv preprint arXiv:2401.11923 (2024).

[173] T. Zhang, H. Feng, W. Chen, Z. Chen, W. Zheng, X. Luo, W. Huang, A. Tung, Chartnavigator: An interactive pattern identification and annotation framework for charts, IEEE Transactions on Knowledge and Data Engineering 35 (2) (2023) 1258–1269.

[174] B. Saleh, M. Dontcheva, A. Hertzmann, Z. Liu, Learning style similarity for searching infographics (2015). `arXiv:1505.01214`.

[175] Y. Ye, R. Huang, W. Zeng, VISAtlas: An image-based exploration and query system for large visualization collections via neural image embedding, IEEE Transactions on Visualization and Computer Graphics (2022) 1–15.

[176] A. Fan, Y. Ma, M. Mancenido, R. Maciejewski, Annotating line charts for addressing deception, in: Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems, Association for Computing Machinery, 2022.

[177] K. Davila, S. Setlur, D. S. Doermann, B. U. Kota, V. Govindaraju, Chart mining: A survey of methods for automated chart analysis, IEEE Transsactions on Pattern Analysis and Machine Intelligence 43 (11) (2021) 3799–3819.

[178] D. J. L. Lee, J. Lee, T. Siddiqui, J. Kim, K. Karahalios, A. G. Parameswaran, You can't always sketch what you want: Understanding sensemaking in visual query systems, IEEE Transactions on Visualization and Computer Graphics 26 (1) (2020) 1267–1277.

[179] S. Cohen, C. Li, J. Yang, C. Yu, Computational journalism: A call to arms to database researchers, in: Fifth Biennial Conference on Innovative Data Systems Research, CIDR 2011, Asilomar, CA, USA, January 9-12, 2011, Online Proceedings, www.cidrdb.org, 2011, pp. 148–151.

[180] Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, B. Guo, Swin transformer: Hierarchical vision transformer using shifted windows, in: Proceedings of the IEEE/CVF International Conference on Computer Vision, 2021, pp. 10012–10022.

[181] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, et al., Learning transferable visual models from natural language supervision, in: International Conference on Machine Learning, PMLR, 2021, pp. 8748–8763.

[182] F.-Y. Sun, J. Hoffman, V. Verma, J. Tang, InfoGraph: Unsupervised and semi-supervised graph-level representation learning via mutual

information maximization, in: International Conference on Learning Representations, 2019.

[183] P. Lewis, E. Perez, A. Piktus, F. Petroni, V. Karpukhin, N. Goyal, H. Küttler, M. Lewis, W.-t. Yih, T. Rocktäschel, et al., Retrieval-augmented generation for knowledge-intensive nlp tasks, in: Proceedings of the International Conference on Neural Information Processing Systems, Vol. 33, 2020, pp. 9459–9474.

[184] J. Liu, J. Jin, Z. Wang, J. Cheng, Z. Dou, J.-R. Wen, RETA-LLM: A retrieval-augmented large language model toolkit, arXiv preprint arXiv:2306.05212 (2023).

[185] A. Baldrati, M. Bertini, T. Uricchio, A. Del Bimbo, Effective conditioned and composed image retrieval combining clip-based features, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2022, pp. 21466–21474.

[186] Y. Ye, Q. Zhu, S. Xiao, K. Zhang, W. Zeng, The contemporary art of image search: Iterative user intent expansion via vision-language model, arXiv preprint arXiv:2312.01656 (2023).

[187] X. Zeng, Z. Gao, Y. Ye, W. Zeng, IntentTuner: An interactive framework for integrating human intents in fine-tuning text-to-image generative models, arXiv preprint arXiv:2401.15559 (2024).

[188] S. E. Kahou, V. Michalski, A. Atkinson, Á. Kádár, A. Trischler, Y. Bengio, FigureQA: An annotated figure dataset for visual reasoning, arXiv preprint arXiv:1710.07300 (2017).

[189] J. Zou, G. Wu, T. Xue, Q. Wu, An affinity-driven relation network for figure question answering, in: IEEE International Conference on Multimedia and Expo, IEEE, 2020, pp. 1–6.

[190] D. H. Kim, E. Hoque, M. Agrawala, Answering questions about charts and generating visual explanations, in: Proceedings of CHI Conference on Human Factors in Computing Systems, 2020, pp. 1–13.

[191] S. Song, J. Chen, C. Li, C. Wang, GVQA: Learning to answer questions about graphs with visualizations via knowledge base, in: Proceedings of the CHI Conference on Human Factors in Computing Systems, 2023.

[192] E. Hoque, P. Kavehzadeh, A. Masry, Chart question answering: State of the art and future directions, Computer Graphics Forum 41 (3) (2022) 555–572.

[193] Z. Yang, X. He, J. Gao, L. Deng, A. Smola, Stacked attention networks for image question answering, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2016, pp. 21–29.

[194] J. Luo, Z. Li, J. Wang, C.-Y. Lin, ChartOCR: Data extraction from charts images via a deep hybrid framework, in: Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision, 2021, pp. 1917–1925.

[195] G. Kim, T. Hong, M. Yim, J. Nam, J. Park, J. Yim, W. Hwang, S. Yun, D. Han, S. Park, OCR-Free document understanding transformer, in: European Conference on Computer Vision, Springer Nature Switzerland, Cham, 2022, pp. 498–517.

[196] J. Cho, J. Lei, H. Tan, M. Bansal, Unifying vision-and-language tasks via text generation, in: International Conference on Machine Learning, PMLR, 2021, pp. 1931–1942.

[197] J. Herzig, P. K. Nowak, T. Mueller, F. Piccinno, J. Eisenschlos, TaPas: Weakly supervised table parsing via pre-training, in: Proceedings of the Annual Meeting of the Association for Computational Linguistics, 2020, pp. 4320–4333.

[198] J. Devlin, M.-W. Chang, K. Lee, K. Toutanova, BERT: Pre-training of deep bidirectional transformers for language understanding, arXiv preprint arXiv:1810.04805 (2018).

[199] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, et al., An image is worth 16x16 words: Transformers for image recognition at scale, arXiv preprint arXiv:2010.11929 (2020).

[200] S. Li, N. Tajbakhsh, SciGraphQA: A large-scale synthetic multi-turn question-answering dataset for scientific graphs, arXiv preprint arXiv:2308.03349 (2023).

[201] C. Andrews, A. Endert, B. Yost, C. North, Information visualization on large, high-resolution displays: Issues, challenges, and opportunities, Information Visualization 10 (4) (2011) 341–355.

[202] J. Walny, B. Lee, P. Johns, N. H. Riche, S. Carpendale, Understanding pen and touch interaction for data exploration on interactive whiteboards, IEEE Transactions on Visualization and Computer Graphics 18 (12) (2012) 2779–2788.

[203] S. K. Badam, F. Amini, N. Elmqvist, P. Irani, Supporting visual exploration for multiple users in large display environments, in: IEEE Conference on Visual Analytics Science and Technology, IEEE, 2016, pp. 1–10.

[204] E. Hoque, V. Setlur, M. Tory, I. Dykeman, Applying pragmatics principles for interaction with visual analytics, IEEE Transactions on Visualization and Computer Graphics 24 (1) (2017) 309–318.

[205] A. Srinivasan, B. Lee, N. Henry Riche, S. M. Drucker, K. Hinckley, InChorus: Designing consistent multimodal interactions for data visualization on tablet devices, in: Proceedings of the CHI Conference on Human Factors in Computing Systems, 2020, pp. 1–13.

[206] J. Tang, Y. Luo, M. Ouzzani, G. Li, H. Chen, Sevi: Speech-to-visualization through neural machine translation, in: SIGMOD Conference, ACM, 2022, pp. 2353–2356.

[207] Y. Lin, H. Li, L. Yang, A. Wu, H. Qu, InkSight: Leveraging sketch interaction for documenting chart findings in computational notebooks, arXiv preprint arXiv:2307.07922 (2023).

[208] A. Kirillov, E. Mintun, N. Ravi, H. Mao, C. Rolland, L. Gustafson, T. Xiao, S. Whitehead, A. C. Berg, W.-Y. Lo, et al., Segment anything, arXiv preprint arXiv:2304.02643 (2023).

[209] S. Liu, Z. Zeng, T. Ren, F. Li, H. Zhang, J. Yang, C. Li, J. Yang, H. Su, J. Zhu, et al., Grounding dino: Marrying dino with grounded pre-training for open-set object detection, arXiv preprint arXiv:2303.05499 (2023).

[210] H. Liu, C. Li, Q. Wu, Y. J. Lee, Visual instruction tuning, arXiv preprint arXiv:2304.08485 (2023).

[211] S. Liu, H. Miao, Z. Li, M. Olson, V. Pascucci, P.-T. Bremer, AVA: Towards autonomous visualization agents through visual perception-driven decision-making, arXiv preprint arXiv:2312.04494 (2023).

[212] J. Lu, B. Pan, J. Chen, Y. Feng, J. Hu, Y. Peng, W. Chen, AgentLens: Visual analysis for agent behaviors in llm-based autonomous systems, arXiv preprint arXiv:2402.08995 (2024).

[213] P. Maddigan, T. Susnjak, Chat2vis: Fine-tuning data visualisations using multilingual natural language text and pre-trained large language models (2023). `arXiv:2303.14292`.

[214] R. Yen, J. Zhu, S. Suh, H. Xia, J. Zhao, Coladder: Supporting programmers with hierarchical code generation in multi-level abstraction, arXiv preprint arXiv:2310.08699 (2023).

[215] A. Agrawal, I. Kajic, E. Bugliarello, E. Davoodi, A. Gergely, P. Blunsom, A. Nematzadeh, Reassessing evaluation practices in visual question answering: A case study on out-of-distribution generalization, in: Findings of the Association for Computational Linguistics: EACL 2023, 2023, pp. 1171–1196.

[216] Z. Ji, N. Lee, R. Frieske, T. Yu, D. Su, Y. Xu, E. Ishii, Y. J. Bang, A. Madotto, P. Fung, Survey of hallucination in natural language generation, ACM Computing Surveys 55 (12) (2023) 1–38.

[217] K. Hu, S. Gaikwad, M. Hulsebos, M. A. Bakker, E. Zgraggen, C. Hidalgo, T. Kraska, G. Li, A. Satyanarayan, Ç. Demiralp, VizNet: Towards a large-scale visualization learning and benchmarking repository, in: Proceedings of the CHI Conference on Human Factors in Computing Systems, 2019, pp. 1–12.

[218] X. Chen, W. Zeng, Y. Lin, H. M. Ai-Maneea, J. Roberts, R. Chang, Composition and configuration patterns in multiple-view visualizations, IEEE Transactions on Visualization and Computer Graphics 27 (2) (2020) 1514–1524.

[219] L. Battle, P. Duan, Z. Miranda, D. Mukusheva, R. Chang, M. Stonebraker, Beagle: Automated extraction and interpretation of visualizations from the web, in: Proceedings of the CHI Conference on Human Factors in Computing Systems, 2018, pp. 1–8.

[220] Y. Ye, R. Huang, W. Zeng, VISAtlas: An image-based exploration and query system for large visualization collections via neural image embedding, IEEE Transactions on Visualization and Computer Graphics (2022) 1–15.

[221] Y. Luo, J. Tang, G. Li, nvbench: A large-scale synthesized dataset for cross-domain natural language to visualization task, arXiv preprint arXiv:2112.12926 (2021).

[222] J. Wang, S. Liu, W. Zhang, Visual analytics for machine learning: A data perspective survey, arXiv preprint arXiv:2307.07712 (2023).

[223] A. I. Anik, A. Bunt, Data-centric explanations: explaining training data of machine learning systems to promote transparency, in: Proceedings of the CHI Conference on Human Factors in Computing Systems, 2021, pp. 1–13.

[224] P. Galanter, Generative art theory, A Companion to Digital Art (2016) 146–180.

[225] H. Wang, P. K. A. Vasu, F. Faghri, R. Vemulapalli, M. Farajtabar, S. Mehta, M. Rastegari, O. Tuzel, H. Pouransari, SAM-CLIP: Merging vision foundation models towards semantic and spatial understanding, arXiv preprint arXiv:2310.15308 (2023).

[226] Y. Yuan, W. Li, J. Liu, D. Tang, X. Luo, C. Qin, L. Zhang, J. Zhu, Osprey: Pixel understanding with visual instruction tuning, arXiv preprint arXiv:2312.10032 (2023).

[227] L. Zhang, A. Rao, M. Agrawala, Adding conditional control to text-to-image diffusion models, in: Proceedings of the IEEE/CVF International Conference on Computer Vision, 2023, pp. 3836–3847.

[228] G. Zheng, X. Zhou, X. Li, Z. Qi, Y. Shan, X. Li, LayoutDiffusion: Controllable diffusion model for layout-to-image generation, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2023, pp. 22490–22499.

[229] C. Liu, Y. Guo, X. Yuan, AutoTitle: An interactive title generator for visualizations, IEEE Transactions on Visualization and Computer Graphics (2023) 1–12.