# Chartist: Task-driven Eye Movement Control for Chart Reading

Danqing Shi
Aalto University
Helsinki, Finland

Yao Wang
University of Stuttgart
Stuttgart, Germany

Yunpeng Bai
National University of Singapore
Singapore, Singapore

Andreas Bulling
University of Stuttgart
Stuttgart, Germany

Antti Oulasvirta
Aalto University
Helsinki, Finland

## ABSTRACT

To design data visualizations that are easy to comprehend, we need to understand how people with different interests read them. Computational models of predicting scanpaths on charts could complement empirical studies by offering estimates of user performance inexpensively; however, previous models have been limited to gaze patterns and overlooked the effects of tasks. Here, we contribute Chartist, a computational model that simulates how users move their eyes to extract information from the chart in order to perform analysis tasks, including value retrieval, filtering, and finding extremes. The novel contribution lies in a two-level hierarchical control architecture. At the high level, the model uses LLMs to comprehend the information gained so far and applies this representation to select a goal for the lower-level controllers, which, in turn, move the eyes in accordance with a sampling policy learned via reinforcement learning. The model is capable of predicting human-like task-driven scanpaths across various tasks. It can be applied in fields such as explainable AI, visualization design evaluation, and optimization. While it displays limitations in terms of generalizability and accuracy, it takes modeling in a promising direction, toward understanding human behaviors in interacting with charts.

## CCS CONCEPTS

• **Human-centered computing** → **HCI theory, concepts and models**; **Information visualization**.

## KEYWORDS

User model; Simulation; Scanpath; Reinforcement learning; LLMs

## 1 INTRODUCTION

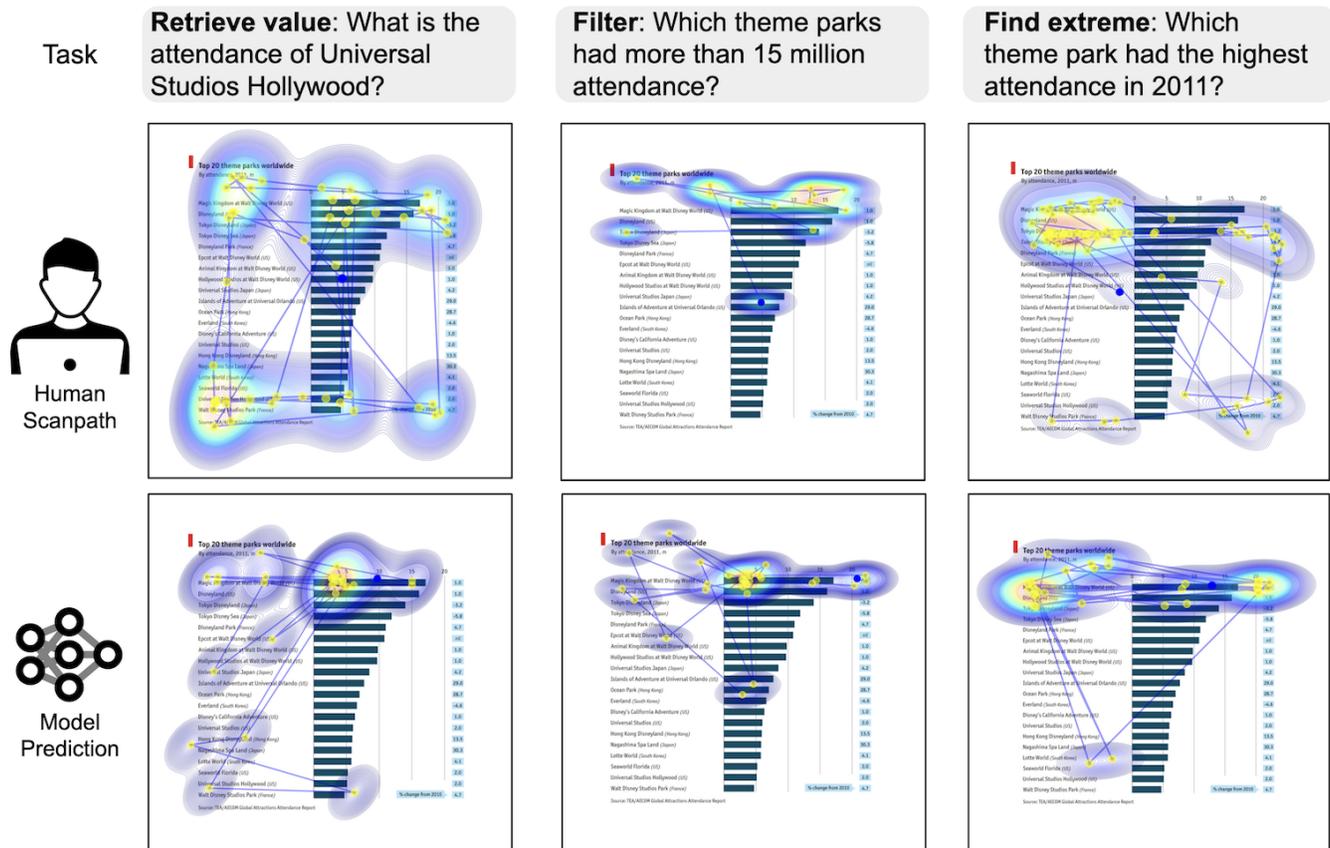Visual attention plays a pivotal role in the field of information visualization [11, 31]. By understanding the visual attention on charts, designers can iteratively refine visualizations using visual saliency as feedback [67]; engineers can enhance their AI models by incorporating human-like attention [69]; and researchers can better understand the link between comprehension and gaze behavior when people read charts [77]. Eye tracking has long been used to understand human visual attention on charts [27]. Beyond visual saliency [66], analyzing human scanpaths provides details of the sequence of fixations, helping researchers understand strategies for reasoning [76]. Previous literature has demonstrated that users observe completely different visual elements when performing different analytical tasks [56]. However, collecting human scanpaths by using eye trackers is expensive both time-wise and monetarily. Simulations are effective in developing theories by rigorously testing user interactions with visual elements in controlled settings. By simulating eye movements, researchers can uncover the mechanisms behind behaviors of users interpreting data visualizations. This enhances understanding of chart reading [51].

Human visual attention is guided by two processes: bottom-up and top-down processes [36]. These processes apply to chart reading as well [83]. Bottom-up attention is driven by salient visual stimuli (e.g., high-contrast colors), whereas top-down attention is task-driven, with specific goals or intentions shaping where and how users focus their attention. However, most visual attention models applied for information visualizations capture only bottom-up (free-viewing) attention [49, 66, 76], thus leaving a gap in understanding how tasks influence human visual attention [4]. Compared to exploratory free viewing, scanpaths for the same analysis task are more coherent; also, they vary greatly between tasks [56].

While recent research has been able to predict task-driven saliency on charts [80], it has remained one step away from addressing how people read charts. Temporal information and individual-specific behaviors are still missing from task-driven-saliency maps. In other words, what would the scanpath look like when a person carries out a particular task on a given chart? In this paper, we present Chartist, the first computational model for predicting task-driven scanpaths on charts [1]. When given both an image of a chart and a sentence as the analytical task, Chartist can simulate a sequence of human-like fixation positions related to the task (see Figure 1). The model has two key distinctions from preexisting models for scanpath prediction: 1) Our model focuses on predicting fixations made during analytical tasks, including both fixation positions and their order, in contrast against the current state-of-the-art models, which concentrate on free-viewing conditions. Task factors' influence makes it challenging to predict task-driven scanpaths via prior models. 2) Unlike goal-driven scanpath predictions, which are
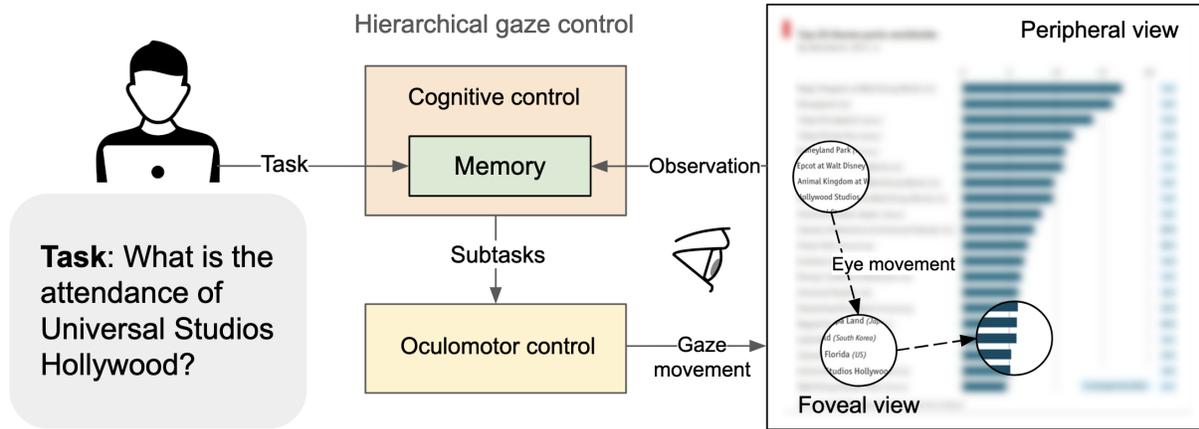
---

[1] https://chart-reading.github.io

**Figure 1: We present CHARTIST, a computational model that can predict task-driven human scanpaths on charts. The figure demonstrates three analytical tasks involved in the study: retrieve value, filter, and find extreme. The visualization illustrates how models' predictions vary across tasks and match the pattern of human scanpaths, with fixation density maps overlaid.**

limited to visual search tasks with natural images [50], analytical tasks require high-level reasoning and also tackling the limitations of human cognition in chart question answering.

To enable such capabilities, we propose a hierarchical control architecture for modeling (see Figure 2). We adopted this architecture for two key reasons: First, it manifests the critical benefit of mirroring how humans break complex tasks into subtasks, which are easier to solve. Studies in cognitive science suggest that humans use hierarchical frameworks in decision-making [13, 26]. Second, the machine-learning community's promising advances in implementing hierarchical architectures in various motor-control domains [15, 35] points to its potential for oculomotor control. This particular hierarchical architecture is composed of a high-level cognitive controller for deciding subtasks and a low-level oculomotor controller for moving the gaze. The cognitive controller is powered by large language models (LLMs) [1] for understanding the task, analyzing the information obtained, and selecting operations for collecting information from the chart. In the oculomotor control, the approach employs deep reinforcement learning (RL) [59], training RL agents for each operation to perform detail-level gaze movements.

We evaluated the model's performance through experiments using human data. We compared scanpath-level similarity with baseline models' output and human scanpaths, where the baselines were the latest general scanpath prediction model [41], scanpath prediction in visual question answering [19], and free-viewing scanpath prediction on charts [76]. The results suggest that the hierarchical gaze control model demonstrates closest similarity to human data across tasks. We also analyze the summary statistic of the gaze behavior from model predictions, which effectively reproduce human-like gaze movement behavior. The evaluation results highlight the potential for the model to exhibit human-like behavior in task-driven scanpaths across different tasks.

In summary, the main contribution of this work is the first computational model CHARTIST, to the best of our knowledge, for predicting task-driven scanpaths on charts. The key technical contribution of CHARTIST is its hierarchical gaze control with a cognitive controller and oculomotor controllers. This architecture enables training the model without relying on human eye movement data. We analyzed human scanpath data from charts, to support modeling for three common analysis tasks: "retrieve value," "filter," and "find extreme". We conducted comprehensive experiments to evaluate scanpath prediction across analytical tasks, comparing scanpath

**Figure 2: The figure illustrates the concept of the model for task-driven eye movement control. When given a task, the agent makes decisions about the next subtask, based on information gathered from observing the chart stored in its memory. Each subtask controls eye movements at pixel level and retrieves information from the foveal vision area of the gaze.**

similarity and providing statistical summaries. Our model performs similarly to humans and better than the baselines in predicting task-driven scanpaths. This study focused on analytical tasks involving statistical charts, where the scanpaths align more with the problem-solving reasoning process and are less influenced by the complexity of the visual representation. At the end of the paper, we discuss the generalizability of the modeling approach.

The paper is structured such that Section 2 reviews prior research into eye tracking connected with information visualizations, analytical tasks, and scanpath prediction models with Section 3 laying further groundwork by introducing the problem formulation and the design of the computational model for predicting task-driven scanpaths, including the hierarchical control architecture and training workflow. We then present our evaluation of the model's performance relative to baseline methods and human data in Section 4. Finally, Section 5 discusses potential applications of the approach and its generalizability.

## 2 RELATED WORK

This section reviews the literature on eye tracking for information visualization, tasks in that domain, and preexisting techniques for scanpath prediction.

### 2.1 Eye Tracking for Information Visualizations

Eye tracking often serves as a proxy for visual perception in human analysis of information visualizations [66]. There is a long history of applying eye tracking techniques to investigate how people perceive visualizations [11, 33, 42, 56]. For instance, Huang [33] examined the relations between eye movement events and visual components in node–link diagrams. Lallé et al. [42] analyzed the connection between gaze behavior and narrative visualizations, and Borkin et al. [11] found links between gaze behaviors over visual elements and the memorability of visualizations. For memorable visualizations, a quick look can already effectively convey the visualization's message. Work by Polatsek et al. [56], demonstrating significant differences in gaze behaviors under three visual analysis

tasks' conditions, attests well that people read completely different regions of charts when handling different tasks. Other scholars have proposed eye-tracking-based approaches for improvements in visual analytics work such as word-sized visualizations [9] and interactive visualizations [53]. However, while eye trackers have become cheaper and more readily available, the scale of eye tracking datasets is still limited because of the amounts of time and money needed [66]. To address this limitation, researchers turned to crowdsourcing platforms as a faster, inexpensive alternative to eye tracking: web cameras [66] and mouse clicks [80] have become popular ways to collect human attention data online. These online alternatives sacrifice the quality of eye tracking data to gain quantity, leaving an open question of how to acquire both high quality and good scale of eye tracking data from information visualizations.

### 2.2 Analytical Tasks in Visualizations

Evidence exists that the task strongly influences how people design and explore visualizations [56, 62, 81]. Amar et al. [4] identified 10 low-level analytical tasks (e.g., retrieving values and finding extremes) while another study [32] highlighted abstract tasks such as background understanding, planning of analysis, and data exploration. Considering specific tasks is critical since visualizations can be created for handling any of the tasks in light of the input data, and also they can be evaluated in terms of how well certain tasks can be accomplished [60]. To facilitate tasks related to visualizations, some researchers have designed visualizations for explicit displaying of data facts [64, 70] – such as showing trends or comparisons directly [61, 63]. Talk2data [30] and Datamator [29] organize data facts related to specific tasks to facilitate question answering. For this study, we adopted three analytical tasks from previous research [56] and conducted further analysis to understand how humans read charts for given tasks. Our analysis inspired us to develop the computational model for predicting task-driven scanpaths over charts.

## 2.3 Scanpath Prediction

In aims of predicting people's spatial and temporal viewing patterns upon exposure to certain stimuli, represented by a sequence of fixations, scholars have studied scanpaths with numerous stimulus types. Among these are natural scenes [20, 82], web pages [22], graphical user interfaces [39], and information visualizations [76]. There is literature on scanpath prediction dedicated to sampling fixations from saliency maps[10, 14, 41]. SaltiNet [6] extended these maps to "saliency volumes," from which sample scanpaths were created, while the models have drawn inspiration from cognitively plausible mechanisms, such as inhibition of return [37, 72] or foveal–peripheral saliency [8, 79]. In HMM-based methods, in turn, the prediction either splits an image into several grids and regards each grid as a single state of observation [75] or classifies the fixations into several states [20].

The advent of deep learning brought new architectures into play for predicting scanpaths: PathGAN's developers [5] proposed using a generative adversarial network (GAN) for scanpath prediction, while Gazeformer [50] encoded stimuli by using a natural-language model, then applied transformer-based modeling to predict visual scanpaths in a zero-shot setting. For task-driven scanpath prediction, Yang et al. [82] put inverse reinforcement learning to use to model human scanpaths during visual search, with Chen et al. [19] likewise proposing an RL model to predict the scanpath during visual question answering.

Modeling scanpaths on information visualizations is challenging, given that viewing behaviors vary greatly across viewers [56]. Wang et al. [76] highlighted the poor performance of prior scanpath prediction models with information visualizations, and they tackled this issue by fine-tuning a multi-duration saliency model [25] for the information's graphical presentation and probabilistically sampling fixations from saliency maps. However, their pioneering model for scanpath prediction in the visualization domain still cannot cope with task-specific scanpaths. To fill the gap, we sought a task-driven model specifically designed to predict a human scanpath over information visualizations.

## 3 CHARTIST: MODELING TASK-DRIVEN EYE MOVEMENT ON CHARTS

This section introduces the problem formulation and presents the computational model of eye movement control on charts in settings of analytical tasks.

### 3.1 Problem Formulation

Given a chart image $C$ and an associated analytical task $x$ stated as text, the model is expected to generate a sequence of fixation positions $\{p_1, p_2, \ldots, p_t\}$. The objective of the output sequence is to closely match the scanpath from humans reading the chart. Specifically, the sequence of fixations represents the visual reasoning process, and the information in the patches of pixels fixated upon should be able to support $x$. We consider general analytical tasks in information visualization [4], and select three of them used in a human eye-tracking data collection [56]:

1) *Retrieve value (RV)*: Given a specific target, find the data value of the target (e.g., what is the value for a certain category?)

2) *Filter (F)*: Given a concrete condition, find which data point satisfies it (e.g., which category has the specific value stated?)
3) *Find extreme (FE)*: Find the data point showing an extreme value for a given attribute within the set of data (e.g., which category shows the highest/lowest value?)
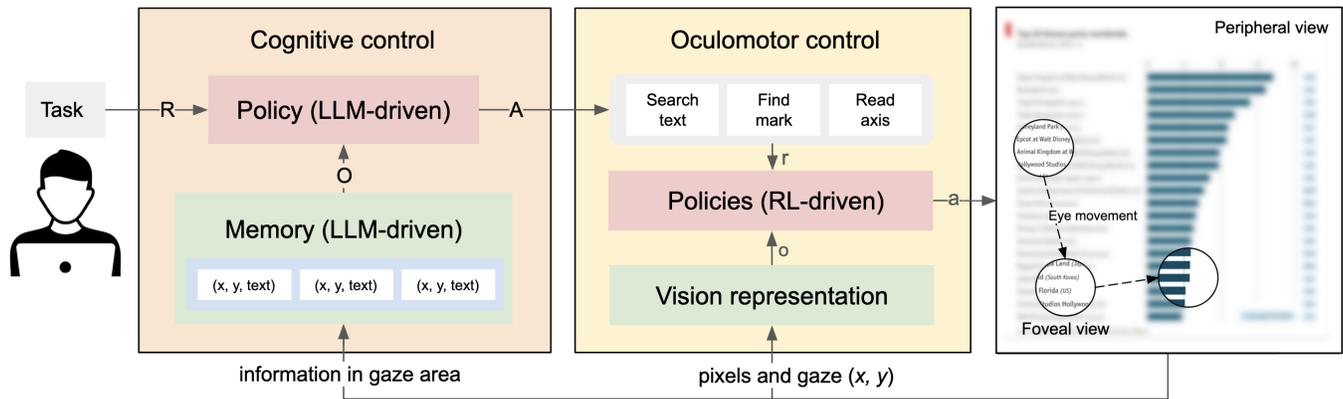
### 3.2 Modeling Overview

Our goal was to develop the model CHARTIST to handle tasks articulated as free-form text and be able to perform gaze movement at a detailed pixel level. We conceptualize the design of the hierarchical gaze control model in Figure 3, where the high-level (cognitive) controller is responsible for reasoning while the low-level (oculomotor) controller determines details of gaze movement. The idea behind this is hierarchical supervisory control [24], which refers to a tiered control system in which the superior controller set goals for its subordinates. The actions from subordinates are integrated into an overall pattern for high-level control [55]. The concept also follows the modeling principle of computational rationality, where we assume that the controllers optimize their policy to maximize expected utility within relevant cognitive bounds [17, 54]. Specifically, the high-level controller handles abstract information processing, comprehension, and memory storage. It sets subtasks to the low-level controller, which then moves the gaze to gather information for task completion. Subsequently, the high-level controller utilizes the amassed information to answer the question.

### 3.3 Cognitive Control

The high-level controller provides cognitive control over the mental processes for a chart, control that performs reasoning in working memory [46]. When performing vision tasks, one observes and analyzes visual information interactively [18]. Throughout this process, people analyze the information in their memory and try to gather more useful information to reduce uncertainty in solving the task. To represent this decision problem accurately, we formulate it as a bounded optimality problem in a partially observable Markov decision process (POMDP). Instead of having access to a full state ($S$) with pixels of the chart associated with the given task, the POMDP expresses a subset of ($S$) as the observation of the model:

- Observation $O$ refers to the information in memory that is captured from eye movements over the chart.
- Action $A$ includes subtasks that the model gives to oculomotor control for performing eye movements.
- Reward $R$ is the correctness of the answer for the task from the chart question answering.

To solve this POMDP, our model uses LLMs for the policy. The rationale behind this choice is that LLMs are well suited to processing higher-level information, as they have been pre-trained on human text data encompassing a wealth of logic related to planning, reasoning, and interaction [34, 43, 74]. Although LLMs are limited in their ability to control low-level motor functions in a precise manner [21], they are proficient at planning and reasoning, with LLaMA [73] and GPT [1] showing impressive language interpretation and reasoning capabilities. Also, recent work has shown that utilizing LLMs in the high-level controllers in hierarchical architecture can produce promising results [15, 34, 45]. For our setting, we used GPT-4o [1] for the policy, which takes the information

**Figure 3: An overview of the hierarchical eye-movement control architecture. When presented with a chart and a task, a cognitive controller, powered by large language models, makes decisions on what to look at next and judges whether it is confident enough to provide an answer to the task's question. It relies on internal memory, which summarizes the information gathered from the chart through eye movements. Once cognitive control has determined the next action, the oculomotor controller is responsible for moving the gaze and observing the chart through a limited vision field. The model's objective is to accurately address the task as quickly as possible within set cognitive and physical constraints.**

accumulated in the memory as the observation and sets subtasks to guide eye movements in order to obtain information needed for solving the task efficiently.

We consider two human limitations when constructing the model's observation: a limited field of vision [23] and memory capacity [47]. The model gets information from the gaze position purely by mimicking the human vision system. An optical character recognition technique [68] is used to extract text from the pixels of the chart, and the text in the gaze area, with the position, is passed to the memory. As a result, the observation consists of image patches (in a limited number) from the full set of chart pixels. The reliability of items in memory is determined by their visit history [44], with overall memory capacity being restricted too. When new information is added to the memory, a previously added item is removed on the basis of a forgetting probability. The probability of forgetting an item in the memory is calculated by means of the formula Softmax($\rho \cdot (t - t_i)$), where $t$ is the current fixation index, $t_i$ is the index of the $i$th item in the memory, and $\rho$ is the weight parameter (set to 0.1 here). The observation is designed as a prompt that summarizes the memory in line with the memory model and explains the model's goal.

Given the summary of the memory information, the LLM policy selects predefined operations for task solving [15, 45]. The operations here are based on a sequence of cognitive stages for charts [27] – 1) *search for text label*: visually searching for a text label or value label related to the task, 2) *find associated mark*: visually searching for a graphical mark of the data point when given a reference label, 3) *read associated value*: visually searching to read the given mark's associated value or textual label. All these actions are allowed to be reused in the process, which enables the model to revisit previous positions for confirmation of the information. Ultimately, if the information in the memory is sufficient to address the task, the gaze movement can stop and an answer can be given. Operations other

than answering the question will be performed by the oculomotor controller for detailed gaze movement.
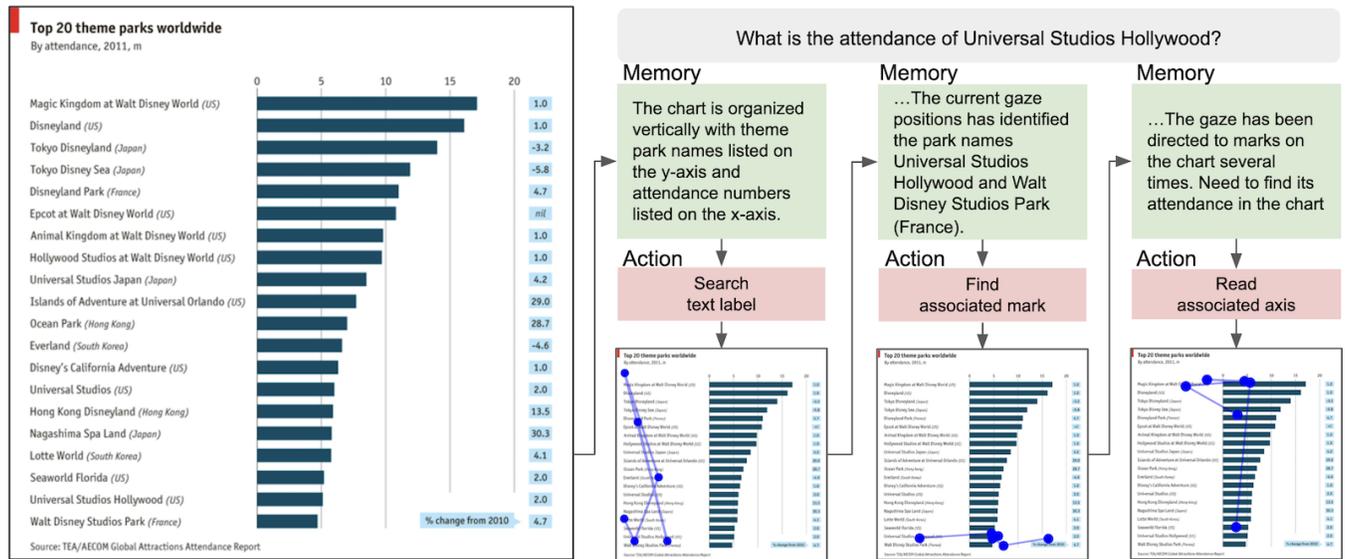
The examples in Figure 4 demonstrate how utilizing memory information and predefined operations aids in scanpath prediction. Model memory uses the summarization capability of LLMs to convert the text and positions gathered to a paragraph as the observation (as shown in the green boxes). The LLM policy then makes decisions and issues subtasks as actions (in red boxes) for the oculomotor control, which performs pixel-level gaze movements.

### 3.4 Oculomotor Control

The oculomotor controller acts as the interface between the cognitive controller and the actual chart-pixel images. Its main function is to control the movement of the gaze over the pixels in order to gather information related to the task at hand. Generating oculomotor behavior at pixel level is another sequential decision-making problem that can be formulated as a POMDP:

- Observation $o$ comprises vision information obtained from the external environment, which is jointly represented by the human vision system and visual short-term memory (VSTM).
- Action $a$ involves specifying the coordinates $(x, y)$ of a particular position to move to.
- Reward $r$ is designed to encourage the gaze to reach the target with less cost. It takes into account the number of target hits as well as the cost associated with the distance of the gaze movement.

Our modeling of a chart reader's observation follows an idea similar to that in visual search [82]. Utilizing a representation for accumulating information through fixations, this employs four components: 1) The foveal and peripheral view come from the human vision system, which receives high-resolution visual input only from the region of the image around the fixation location. It includes two pixel-based modules to read the chart: foveal and

**Figure 4: The figure gives examples of how the internal memory helps the cognitive controller to remember what has been read and then select actions for detailed gaze movement. A green box indicates the information held in memory, a red box represents the action selected by cognitive control, and the blue lines in the images reflect the eye movement scanpaths.**

peripheral vision [23]). 2) Visual saliency provides a bottom-up signal to a chart reader for the given task. The saliency of the chart affects gaze behavior. We use a task-driven saliency model to represent this feature [80]. 3) Visit history represents VSTM, which stores visual information for a few seconds, thereby allowing its use in ongoing cognitive tasks [3]. We represent this history through a matrix where each point is marked as visited or not. 4) A goal-related reference position serves as the initial starting point of gaze movement. For example, the reader might begin at the position of a text label for locating the associated graphical mark, where the position of the text label serves as the reference for the sub-goal. We use a one-hot matrix to represent the reference, in which all cell values are 0 apart from the single 1 that identifies the target. All these components are encoded together via the deep convolutional neural network, followed by a fully connected network.

We train reinforcement learning policies to solve the POMDP for the oculomotor control, because it has been proven to effectively address decision-making challenges in prediction of details of gaze movement [7, 38, 65, 82]. In our detail-level implementation, we resize the input chart images to be $320 \times 320$ and discretize the fixation position into a $20 \times 20$ map. Consequently, each fixation becomes a $16 \times 16$ image patch, and the gaze position is randomly sampled from within that patch. In this setup, the maximum approximation error resulting from this discretization process is less than one degree of the visual angle [82]. Ultimately, both the scanpath and the image will be converted back to the original chart size from $320 \times 320$ pixels.

## 3.5 Workflow

Our implementation of CHARTIST is trained and tested on a collection of tasks and charts. There are four steps, illustrated in Figure 5. In Step 1, real-world charts a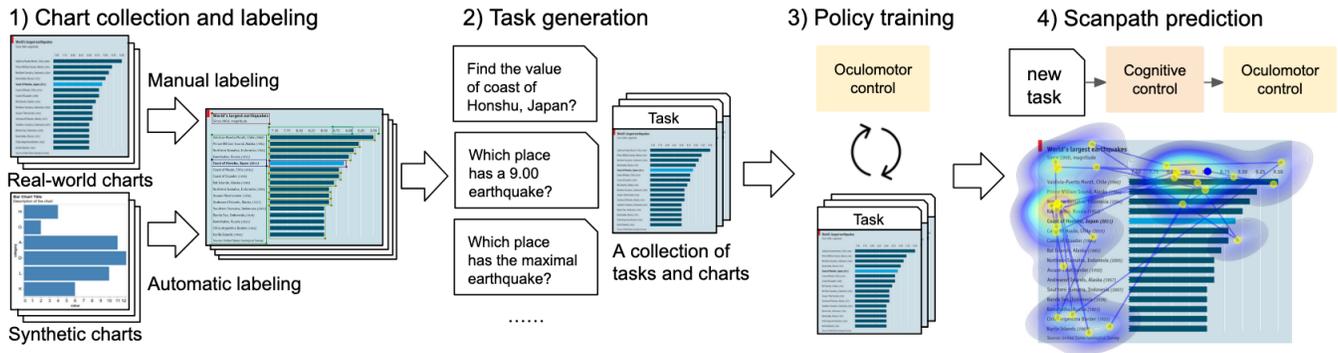re manually collected and labeled for areas of interest (AOIs), while synthetic charts are automatically generated and labeled in a manner powered by Vega-Lite [58]. The inclusion of synthetic charts helps increase the diversity of the chart collection and addresses the challenge of obtaining numerous annotated charts. In Step 2, tasks are automatically generated in line with specific rules for the *RV*, *F*, and *FE* tasks. These tasks and labeled charts constitute a data collection for the training environment.t With Step 3, the policies for oculomotor control are trained through reinforcement learning (using proximal policy pptimization, PPO [59]) to optimize gaze movements, enabling the system to reach task-relevant positions as quickly as possible while adhering to vision constraints. Importantly, no eye tracking data are required for PPO training. In the last phase, prediction, the hierarchical architecture combines pre-trained LLMs (GPT-4o) for cognitive control with RL policies for oculomotor control to generate the scanpath prediction.

## 4 EXPERIMENTS

This section presents the experiments conducted for evaluating and comparing scanpath prediction models. We evaluated CHARTIST specifically in terms of scanpath similarity and statistical summaries of eye movement behavior. The results from evaluation of our model in comparison to the baselines are summarized in Table 1.

## 4.1 Data and Metrics

We evaluated CHARTIST by using 12 distinct analytical tasks with horizontal bar charts from a task-driven scanpath dataset [56]. Each task has at least 14 human scanpaths ($M = 15.25$, $SD = 0.60$), for a total of 183 human scanpaths. This set of ground-truth human data was collected by means of Tobii X2-60 eye trackers at 60 Hz

**Figure 5: An overview of the training workflow: 1) chart collection and labeling, wherein diverse real-world and synthetic charts are gathered, involving manual and automatic annotation of AOIs; 2) task generation, utilizing a rule-based approach to create tasks based on labeled charts to construct a data collection for training; 3) policy training, in which policy models are trained via RL from chart images with tasks; and 4) scanpath prediction, wherein pre-trained LLMs and RL policies are coordinated hierarchically to predict task-driven gaze movements over charts.**

while participants were engaged in these analytical tasks on 24.1-inch monitors at a resolution of 1920 × 1080 pixels. To evaluate model performance, we compared the generated scanpaths against human ground truth, with the aim of ascertaining whether the models can produce human-like scanpaths and replicate natural gaze movement patterns. We should stress that, to ensure unbiased evaluation, none of the human eye movement data informed our training of the models, only their assessment. Furthermore, except for scanpath prediction that relies on human-generated data as the training set, no such eye movement data were involved in training for our approach.

Following recent practice in scanpath prediction [71, 76], we evaluated CHARTIST by using three established scanpath-based metrics: dynamic time warping (DTW), Levenshtein distance (LEV), and Sequence Score. DTW computes the optimal alignment between two scanpaths, where lower values indicate better correspondence. For this paper, DTW was calculated in two-dimensional position coordinates. Both LEV and Sequence Score represent the semantic order of scanpaths as sequences of letters by mapping each fixation to a unique letter, then measuring the string-editing distance between the sequences [52]. In LEV, letters are defined by the grid regions on which fixations land. For Sequence Score [76, 82], letters are based on areas of interest (AOIs) such as the title and legend. Sequence Score values are normalized between 0 and 1, with higher scores reflecting better alignment. For DTW, LEV, and Sequence Score, we report the *mean* and *best* evaluation scores (see Table 2). For each method, the number of scanpaths equaled that of human scanpaths. The *mean* scores are the averages across all human–predicted scanpath pairs, while the *best* ones represent the maximum of all pairs for each prediction [19, 76].

We looked beyond scanpath metrics, introducing more detailed measurements inspired by Goldberg and Helfman [28] to show a statistical summary of task-driven scanpaths over charts.

- *Number of fixations*: The length of a scanpath can be measured as the count of gaze fixations (between motions, or saccades). According to the human data, the number of fixations in task-driven scanpaths over charts (89.8 on average)

is much larger than the number in free-viewing tasks (37.4 on average). This reflects the difficulty of analytical tasks relative to free viewing of charts.

- *Fixation on task-dependent AOI ratio*: Task-dependent AOIs are regions that are relevant to the task, such as value labels, text labels, and data points [56]. People's focus on these areas indicates how they are processing the task. Inspired by the Hit Any AOI Rate metric [78], this measurement provides a summary of the overall visual attention to task-related regions. Although the scanpath is task-driven, we ascertained that most of the eye movement does not occur in task-dependent regions: fewer than 20% of fixations fell in task-dependent AOIs. This suggests that people might devote more time to gathering information or confirming it.

- *Percentage of fixations within each area*: We considered the percentage of time devoted to looking at distinct parts of a chart – namely, the key areas of charts: the title, the marks (such as bars or data points), and the axes. This assists in summarizing where people are focusing their visual attention. The percentages are calculated by dividing the number of fixations in a specific area by the total number of fixations. Humans direct most of their fixations to the axes, then the region of graphical marks. This might be because the three analysis tasks probed are strongly related to values, not other visual features.

- *Fixation transitions*: We also used a metric capturing the average number of times the eyes move from one area to another during a task. It helps us understand how often the eyes' fixations shift between distinct areas. Frequent fixation transitions may point to room for improvement in the design of the chart, such as bringing related elements closer together. From human data, we identified a high number of fixation transitions (about 20 per task). We found that, on average, about four consecutive fixations follow each fixation transition.

- *Revisit frequencies*: The average number of fixations returning to a previously visited area during a task proved similarly

**Table 1: A quantitative benchmark of the task-driven human scanpaths on bar charts. The best results are shown in dark green. All results within 1 standard deviation from human are in light green.**

| Task | Metric | Human | CHARTIST | VQA [19] | UMSS [76] | DG iii [41] |
|---|---|---|---|---|---|---|
| Retrieve value | Sequence Score ↑ (*mean*) | 0.486 | **0.413** | 0.127 | 0.246 | 0.345 |
| | Sequence Score ↑ (*best*) | 0.638 | **0.472** | 0.151 | 0.312 | 0.404 |
| | LEV ↓ (*mean*) | 154.7 | **151.8** | 164.6 | 157.4 | 189.7 |
| | LEV ↓ (*best*) | 121.5 | **141.7** | 162.5 | 150.9 | 181.4 |
| | DTW ↓ (*mean*) | 26,018 | 28,172 | 36,703 | **26,305** | 37,160 |
| | DTW ↓ (*best*) | 17,692 | 23,541 | 33,121 | **21,081** | 30,574 |
| | Number of fixations | 88.6 (57.0) | **48.4** (4.9) | 8.5 (0.5) | 22.2 (3.7) | - |
| | Task AOIs ratio (%) | 10.4 (7.8) | **4.1** (4.3) | 1.2 (3.9) | **3.8** (5.8) | **5.8** (8.2) |
| | Fixation-on-title ratio (%) | 10.5 (9.3) | **5.3** (3.3) | 2.8 (5.7) | **13.7** (8.7) | 33.5 (13.4) |
| | Fixation-on-mark ratio (%) | 31.5 (13.9) | **38.6** (21.3) | 64.0 (24.3) | 51.1 (12.8) | 15.3 (10.3) |
| | Fixation-on-axis ratio (%) | 47.5 (19.0) | **43.9** (17.2) | 23.8 (20.8) | 21.5 (14.7) | **28.6** (10.5) |
| | Fixation transitions | 20.1 (12.4) | **16.7** (7.0) | 3.4 (1.7) | **10.6** (3.3) | 33.9 (9.2) |
| | Revisit frequency title | 2.6 (2.4) | **2.0** (1.3) | **0.2** (0.5) | **2.1** (1.2) | 10.1 (2.6) |
| | Revisit frequency mark | 7.9 (5.1) | **6.8** (3.4) | 1.8 (0.8) | 4.3 (1.5) | **8.6** (4.4) |
| | Revisit frequency axis | 7.9 (4.4) | **7.8** (3.4) | 1.3 (1.0) | 3.0 (1.9) | **12.4** (3.2) |
| Filter | Sequence Score ↑ (*mean*) | 0.452 | **0.379** | 0.149 | 0.271 | 0.321 |
| | Sequence Score ↑ (*best*) | 0.655 | **0.460** | 0.175 | 0.340 | 0.382 |
| | LEV ↓ (*mean*) | 165.8 | **155.3** | 167.6 | 161.0 | 201.5 |
| | LEV ↓ (*best*) | 114.5 | **146.4** | 164.7 | 153.5 | 192.6 |
| | DTW ↓ (*mean*) | 23,732 | 24,617 | 29,141 | **22,817** | 38,139 |
| | DTW ↓ (*best*) | 15,238 | 19,879 | 25,233 | **17,725** | 30,981 |
| | Number of fixations | 89.1 (68.7) | **48.7** (3.3) | 8.5 (0.5) | **22.2** (3.7) | - |
| | Task AOIs ratio (%) | 19.1 (14.4) | **10.6** (11.2) | **13.9** (17.4) | 3.8 (5.8) | **5.8** (8.2) |
| | Fixation-on-title ratio (%) | 5.5 (5.7) | **4.7** (2.3) | **2.0** (5.0) | 13.7 (8.7) | 33.5 (13.4) |
| | Fixation-on-mark ratio (%) | 41.9 (18.4) | **49.1** (21.8) | 68.8 (23.5) | **51.1** (12.8) | 15.3 (10.3) |
| | Fixation-on-axis ratio (%) | 43.3 (20.9) | **36.7** (19.4) | 20.6 (19.4) | 21.5 (14.7) | **28.6** (10.5) |
| | Fixation transitions | 18.4 (13.3) | **15.3** (5.8) | 3.2 (1.7) | **10.6** (3.3) | 33.9 (9.2) |
| | Revisit frequency title | 1.4 (1.7) | **2.0** (1.1) | **0.2** (0.4) | **2.1** (1.2) | 10.1 (2.6) |
| | Revisit frequency mark | 7.7 (5.6) | **6.0** (2.9) | 1.8 (0.8) | 4.3 (1.5) | **8.6** (4.4) |
| | Revisit frequency axis | 8.0 (5.5) | **7.0** (2.5) | 1.2 (1.1) | **3.0** (1.9) | **12.4** (3.2) |
| Find extreme | Sequence Score ↑ (*mean*) | 0.454 | **0.378** | 0.126 | 0.253 | 0.362 |
| | Sequence Score ↑ (*best*) | 0.627 | **0.457** | 0.149 | 0.320 | 0.428 |
| | LEV ↓ (*mean*) | 167.3 | **151.4** | 168.3 | 160.5 | 188.9 |
| | LEV ↓ (*best*) | 123.9 | **143.9** | 168.3 | 155.6 | 180.9 |
| | DTW ↓ (*mean*) | 27,701 | 26,677 | 34,287 | **25,398** | 36,878 |
| | DTW ↓ (*best*) | 18,626 | 22,537 | 31,379 | **20,631** | 30,400 |
| | Number of fixations | 91.9 (64.9) | **46.2** (6.8) | 8.5 (0.5) | 22.2 (3.7) | - |
| | Task AOIs ratio (%) | 1.7 (3.4) | **0.2** (0.7) | 0 (0) | **3.8** (5.8) | 5.8 (8.2) |
| | Fixation-on-title ratio (%) | 12.0 (8.6) | **5.4** (3.2) | 2.4 (5.5) | **13.7** (8.7) | 33.5 (13.4) |
| | Fixation-on-mark ratio (%) | 34.4 (19.4) | **41.6** (15.6) | 64.7 (25.6) | 51.1 (12.8) | **15.3** (10.3) |
| | Fixation-on-axis ratio (%) | 39.7 (22.8) | **40.7** (11.3) | **23.7** (22.9) | 21.5 (14.7) | **28.6** (10.5) |
| | Fixation transitions | 22.6 (16.8) | **16.0** (4.9) | 3.0 (1.6) | **10.6** (3.3) | **33.9** (9.2) |
| | Revisit frequency title | 3.2 (3.3) | **2.0** (1.1) | **0.2** (0.5) | **2.1** (1.2) | 10.1 (2.6) |
| | Revisit frequency mark | 8.1 (5.6) | **6.1** (2.4) | 1.7 (0.8) | 4.3 (1.5) | **8.6** (4.4) |
| | Revisit frequency axis | 8.9 (5.7) | **7.7** (2.3) | 0.8 (1.1) | 3.0 (1.9) | **12.4** (3.2) |

revealing. Human data exhibited high revisit rates. Spatially, users revisit marks and also axes eight times, on average. This frequent double-checking of data information in the chart for the answer may be due to forgetting the information.

These metrics help us evaluate whether the model's predictions can accurately capture general human patterns followed with charts for particular tasks. We strove for a system in which the predicted scanpath closely matches human ground-truth performance, ideally being within one standard deviation of the mean value.

## 4.2 Comparison Methods

Given the lack of existing methods for predicting task-driven scanpaths on information visualizations, we compare Chartist against human ground truth with three closely related baselines:

- *Human* [56]. With the scanpath metrics, we conducted leave-one-out cross-validation among the human scanpaths. For each viewing condition, every human scanpath was compared with all other human scanpaths for similarity. Human scanpaths were compared with themselves for the *mean* scores but not for contributions to the *best* scores. In applying the statistical metrics, we treated all the human data as the ground truth and gauged all modeling methods by their closeness to this ground truth.
- *VQA scanpaths* [19]. VQA is a deep reinforcement learning model that predicts human visual scanpaths in the context of images with visual question answering. The paper reporting on it demonstrates its strong generalizability across various tasks and datasets, indicating optionality as an approach for predicting task-driven scanpaths over charts.
- *UMSS* [76]. UMSS represents the state-of-the-art scanpath prediction model for visualizations, making it the most relevant work in this area. However, it is designed to predict scanpaths in a free-viewing context for information visualizations, rather than consider specific tasks. Its inclusion allows for comparison between scanpath prediction with and without task-linked factors.
- *DeepGaze iii* [41]. DeepGaze iii is a deep-learning-based model that integrates image data with information about previous fixations to forecast free-viewing scanpaths over static images. Trained on large sets of eye tracking data from natural images, it serves as a baseline for evaluating the effect both of stimuli and of tasks on scanpaths.

## 4.3 Results

*Chartist demonstrates high similarity in scanpaths.* The first six rows for each task type in Table 1 present the results from our three scanpath similarity metrics. Chartist achieved the highest performance by the Sequence Score and LEV metrics, and it ranked second for DTW, with scores closely approaching the maximum and also close to human ground truth. Specifically, Chartist closely approximates the latter Sequence Score in terms of mean performance, achieving a score of 0.413, relative to 0.486. The UMSS method, while securing first place for DTW, ranked second for LEV and third for Sequence Score. Chartist outperforms human ground

truth in LEV *mean* (151.8 vs. 154.7). This means that the predictions from Chartist deviate less from the "average human scanpath."

The results from the scanpath metrics show that Chartist performed better than the baselines by the Sequence Score and LEV metrics, which are based on regions, but not the DTW metric, which is based on pixel-wise distances. This suggests that Chartist is more similar to human data when one factors in the semantic order of fixation positions in meaningful portions of charts. However, it may not fully match human data for pixel-level similarity. This result is consistent with the discussion in the literature [76], which has concluded that metrics based on pixel-wise distances between scanpaths might not wholly capture the quality of human scanpaths. Therefore, we must conduct further analysis of the statistical summary of eye movement behaviors.

*Chartist aligns more closely with human statistical patterns.* The last nine rows for each task type in Table 1 provide the mean and standard deviation for each eye movement behavior. Because UMSS and DeepGaze iii are not task-driven, our analysis used the same predicted scanpaths across all tasks. Chartist achieves strong alignment with human data, with all 27 of its values for the eye movement behavior metrics falling within one standard deviation of the human mean and with 18 of them being the closest to the human data's mean. In comparison, 18 of UMSS's 27 values lie within one standard deviation, and five of them are the closest to the mean. The corresponding figures for DeepGaze iii are 13 out of 27 and 4, respectively, while VQA yielded only six values within the range and only one of the 27 was the closest to the human mean.

Examining the detailed metrics across tasks reveals that humans show significantly variable task-dependent AOI ratios. They devote the majority of their fixations to task AOIs when performing the *F* task (19.1%). That is followed by the *RV* task (10.4%), with the *FE* task having the lowest percentage (1.7%). This distribution makes sense: the first two tasks require individuals to focus on a specific data label, while *FE* can be completed by directly observing the general shape of the graph. Chartist is the only model that successfully replicates this phenomenon by reproducing the human order of task-dependent AOI ratios: FE (10.6%), then RV (4.1%), and finally FE (0.2%). As for per-region fixation ratios, humans direct the most fixations to axes, followed by marks, across all three tasks. Chartist successfully reproduces this phenomenon in the case of the *RV* task. For the *F* and *FE* tasks, Chartist shows similar distributions. In contrast, the VQA and UMSS baselines consistently allocate over 50% of fixations to the marks, and DeepGaze iii allocates most fixations to the title. In the realm of revisits, Chartist and DeepGaze iii align with human data, revisiting the axes most frequently, while the other two models revisit the marks most often. In summary, our analysis demonstrates that Chartist exhibits the pattern most similar to human data.

*Qualitative analysis.* Figure 6 showcases predicted scanpaths from Chartist and the three baseline models across six tasks, with fixation density maps overlaid. In all cases, our model's predictions are closer to the human data than the baseline models'. Chartist and VQA scanpaths both are task-driven, unlike UMSS and DeepGaze iii's, which cannot predict scanpaths solely from images. Here are the main observations from Figure 6:
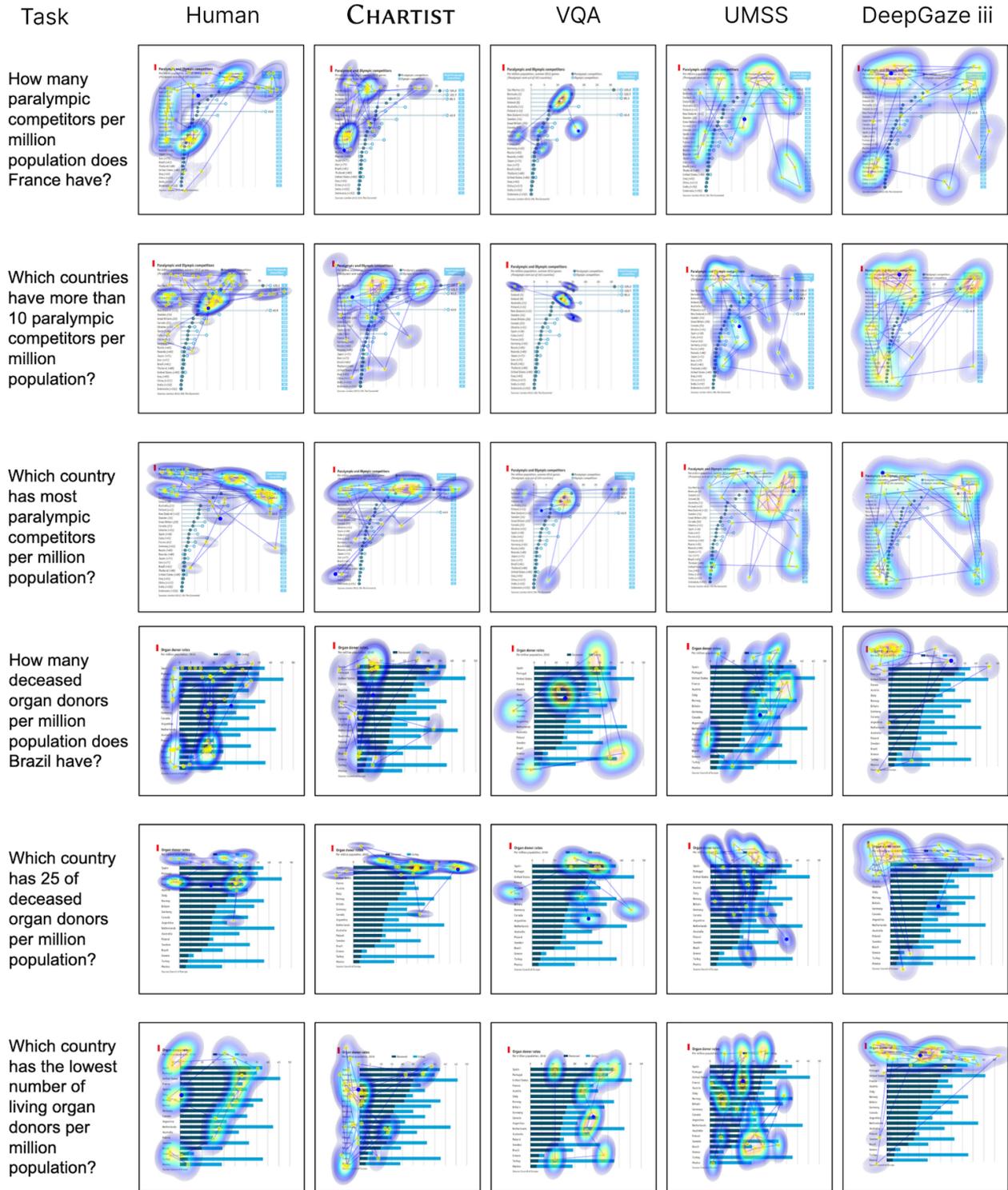
**Figure 6: Qualitative comparison: for three tasks, an illustration of CHARTIST's predictions relative to three baselines – VQA scanpaths [19], UMSS [76], and DeepGaze iii [41]. CHARTIST is able to capture human scanpath patterns displayed during analytical tasks.**

- DeepGaze iii predicts scanpaths from the bottom up on the basis of the saliency of a natural image. That results in a high number of fixations on the title, consistent with producing the highest fixation-on-title ratios.
- Although UMSS is not task-driven, its predicted scanpath shape closely mirrors human data. This indicates that, even in task-driven scenarios, bottom-up mechanisms exert a significant influence.
- Nevertheless, Chartist, as does VQA, captures human gaze patterns more effectively across tasks than these free-viewing models. Importantly, Chartist outperforms the VQA scan-path model, for VQA often predicts fixating on irrelevant areas, in the absence of specific knowledge of visualization structures.

The examples in Figure 1 demonstrate how Chartist's predictions compared to human data for the three tasks. The chart read displays a list of top-ranked theme parks worldwide with their corresponding attendance numbers for 2011. When given the *RV* task of answering "what is the attendance level of Universal Studios Hollywood?" both the human user and the model focus on the text labels to find the theme park in the chart and the relevant positions for the mark and on the value axis. For the *F* task, the human user and the model both frequently look at the value axis. Regarding the *FE* task, both human and model focus on the top of the mark and also fixate on the text label associated with that mark. We noted that human eye movements are also drawn to other text labels, such as annotations and textual descriptions, while Chartist remains task-focused without getting distracted by unrelated information.

## 5 DISCUSSION

While the results show that Chartist is able to simulate human-like eye movements when performing analytical tasks, there is a need to expand on our discussion of the model's practical implications, the generalizability of the modeling approach, and the limitations and potential for supporting sophisticated chart-based question answering.

### 5.1 Applications

*Visualization design evaluation.* Chartist can assist in evaluation of chart design. With well-controlled experiment conditions, eye tracking data afford valuable insight into chart designs, especially relative to alternative designs. For example, Goldberg and Helfman [27] showcased eye tracking's value in comparing line and radial graphs for reading of values, by allowing researchers to understand the viewing order of AOIs and the task completion time. Chartist holds potential to replace human input to evaluation based on eye tracking. With the simulated scanpaths from Chartist, chart designers can obtain quick and cost-effective feedback that yields the benefits from eye tracking without requiring an expensive empirical study.

*Visualization design optimization.* Beyond evaluation, another potential usage application of Chartist is to help optimize visualization design [67]. Like other fields of design, visualization design requires user feedback for continual iteration. When visualization designers create charts for specific tasks, they may wonder if the design is suitable for delivery. With the predicted scanpaths from the model, they can easily access quick and affordable feedback before deeming a candidate design ready for expensive evaluation in a user study. Predictive models could offer feedback to designers or even provide optimization goals in automated visualization design frameworks. The ultimate goal is grounding for recommendations for visualizations that support specific tasks [2] and even automation of visualization design in real time. Today's human-in-the-loop design optimization paradigm [40] could shift to a user-agent-in-the-loop approach, wherein a computational agent that simulates human feedback enables scalable and efficient design evaluation.
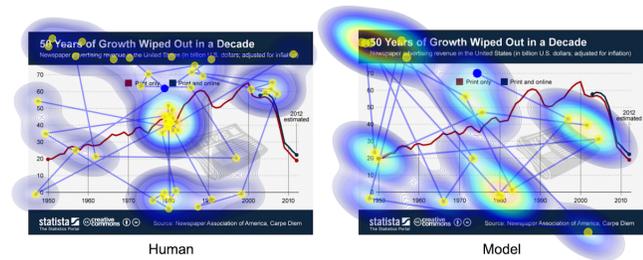
*Explainable AI in chart question answering.* Systems for answering questions via charts [48] are typically viewed as black boxes that generate answers directly from a given chart and natural-language question. In contrast, Chartist introduces a glass-box approach that answers questions through a step-by-step reasoning process. This method enhances the alignment between human and machine attention [69]. We anticipate that this approach could lead to significant improvements in chart question answering [48] and greater compatibility with explainable AI systems.

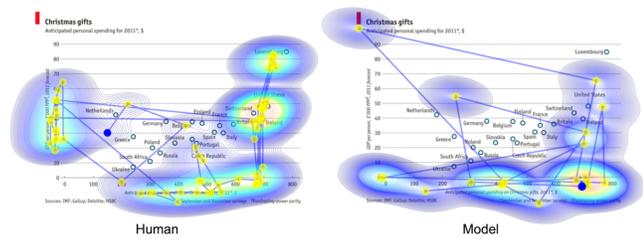### 5.2 Extending the Model beyond Bar Charts

Our modeling approach can be extended to many visualization types besides bar charts. We analyzed the visualization taxonomy outlined in prior work [11, 12], including area, circle, diagram, distribution, grid, line, map, point, table, text, tree, and network, then categorize these visualization techniques into two groups: those that are feasible to extend with minor changes and those that are out of reach, requiring additional features.

Our modeling approach can be applied to most statistical charts either directly or upon rectification of minor issues. For instance, extending the model to interpret *line charts* and *area charts* is feasible when the axis labels are clearly defined. The trend patterns of lines and areas can be perceived by the peripheral vision as visual guidance. For *point charts*, such as scatterplots, the model performs well in conditions of sparse data points. However, individual points may be obscured in dense scatterplots, making it difficult to label data when points are cluttered or overlapping. *Distribution charts*, such as histograms, and *circle charts*, such as pie charts, are similar to bar charts in that they use the area of marks to represent values. Retrieving exact values from these two presentation types can be imprecise on account of the ranges of the bins and inaccuracies in estimating angles or arc lengths. Reading *grid charts* (e.g., heatmaps) too is feasible; however, identifying the values necessitates understanding color intensity, a factor that can sometimes lead to ambiguity. Modeling scanpaths on *tables* or *text* for retrieval tasks is tractable under the current modeling approach, but a lack of visual pattern recognition may render the results poor. To further examine the generalizability of this category, we considered two additional cases, using a line chart and a scatterplot. We manually labeled the charts, trained the model, and made predictions. As Figure 7 attests, the trained model performs well for these two chart types when compared to human ground-truth scanpaths.

Other, sophisticated visualization types are out of reach because they require additional features, particularly prior knowledge and advanced reasoning abilities. For instance, reading *maps* involves associating spatial regions with colors, sizes, or symbols to retrieve

(a) An *RV* task with a line chart: "What was the revenue from newspaper advertising in 1980?"



(b) An *F* task with a scatterplot: "In which countries do people anticipate spending about $700 for personal Christmas gifts?"

**Figure 7: Two cases that illustrate the generalizability of the modeling approach, showing the extension of Chartist to a line chart and a scatterplot. The model's predictions are spatially similar to human ground-truth scanpaths.**

related values. Also, when interpreting maps, people rely heavily on preexisting geographical knowledge as a basis for efficient visual searches. Complex designs with intricate structures, such as *diagrams*, *trees*, and *network graphs*, typically require advanced reasoning based on connections. All these skill requirements point to a need for further study in this area.

### 5.3 Paths toward Sophisticated Tasks in Chart Question Answering

Although the model focuses primarily on gaze prediction, it is worth exploring potential improvements for enriching its sophisticated question answering capabilities. We also discuss its limitations.

Our current model does not achieve the same level of accuracy as the state-of-the-art models represented by the ChartQA benchmark [48]. Unlike other models that can access the full chart image, Chartist is limited by its foveal vision and restricted spatial reasoning abilities. For instance, if a bar's height falls between two labeled values, such as 10 and 15, the model might choose either 10 or 15 as its answer when interpreting the axis, failing to provide a more precise value. This limitation stems from the constrained spatial perception capabilities of LLMs, which are central to cognitive control. One possible solution is integrating multi-modal LLMs [16], for which recent research has demonstrated an accuracy rate of 81.3%.

The sense-making process for complex visualizations may be inherently challenging. Even humans often struggle with understanding how the data are encoded, recognizing a given chart's

purpose, tackling readability issues, performing numerical calculations, identifying relationships among data points, and navigating the spatial arrangement of graphical elements [57]. Our model is designed to be straightforward and objective, focusing on analysis tasks related to statistical charts, but it does not fully capture the complexities of visualizations. A possible enhancement in this respect would be to integrate the model with human sense-making practices [57] or to incorporate a framework of human understanding [2]. Such integration could facilitate better simulation of a human-like problem-solving process.

## 6 CONCLUSION

Our work contributes a computational model Chartist that simulates the eye movements on charts when humans solve visual analytical tasks. The model benefits greatly from its hierarchical gaze control architecture wherein the high-level cognitive controller performs reasoning using memory while the low-level oculomotor controller directs the gaze within visual constraints. By following the principle of computational rationality, we are able to train the model in a controlled environment instead of relying on human eye tracking data. This circumvents the costly and time-consuming process of gathering such data. The results indicate that the model is better than baselines at generating human-like eye movements during analytical tasks. The predicted scanpaths closely match the spatial positions and temporal order of human scanpaths. While there are limitations to its generalizability and accuracy in question answering, it paves the way for further advances in modeling-based approaches.

## REFERENCES

[1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. GPT-4 technical report. *arXiv preprint arXiv:2303.08774* (2023).

[2] Danielle Albers, Michael Correll, and Michael Gleicher. 2014. Task-driven evaluation of aggregation in time series visualization. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*. 551–560.

[3] George A Alvarez and Patrick Cavanagh. 2004. The capacity of visual short-term memory is set both by visual information load and by number of objects. *Psychological Science* 15, 2 (2004), 106–111.

[4] Robert Amar, James Eagan, and John Stasko. 2005. Low-level components of analytic activity in information visualization. In *IEEE Symposium on Information Visualization (INFOVIS)*. 111–117.

[5] Marc Assens, Xavier Giro-i Nieto, Kevin McGuinness, and Noel E O'Connor. 2018. PathGAN: Visual scanpath prediction with generative adversarial networks. In *Proceedings of the European Conference on Computer Vision (ECCV)*. 406–422.

[6] Marc Assens Reina, Xavier Giro-i Nieto, Kevin McGuinness, and Noel E O'Connor. 2017. Saltinet: Scan-path prediction on 360 degree images using saliency volumes.

In *Proceedings of the IEEE International Conference on Computer Vision Workshops*. 2331–2338.

[7] Yunpeng Bai, Aleksi Ikkala, Antti Oulasvirta, Shengdong Zhao, Lucia J Wang, Pengzhi Yang, and Peisen Xu. 2024. Heads-Up Multitasker: Simulating attention switching on optical head-mounted displays. In *CHI '24: Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*. Article 79.

[8] Wentao Bao and Zhenzhong Chen. 2020. Human scanpath prediction based on deep convolutional saccadic model. *Neurocomputing* 404 (2020), 154–164.

[9] Fabian Beck, Tanja Blascheck, Thomas Ertl, and Daniel Weiskopf. 2015. Word-sized eye-tracking visualizations. In *Workshop on Eye Tracking and Visualization (ETVIS)*. 113–128.

[10] Giuseppe Boccignone and Mario Ferraro. 2010. Gaze shifts as dynamical random sampling. In *2010 2nd European Workshop on Visual Information Processing (EUVIP)*. IEEE, 29–34.

[11] Michelle A Borkin, Zoya Bylinskii, Nam Wook Kim, Constance May Bainbridge, Chelsea S Yeh, Daniel Borkin, Hanspeter Pfister, and Aude Oliva. 2015. Beyond memorability: Visualization recognition and recall. *IEEE Transactions on Visualization and Computer Graphics (TVCG)* 22, 1 (2015), 519–528.

[12] Michelle A Borkin, Azalea A Vo, Zoya Bylinskii, Phillip Isola, Shashank Sunkavalli, Aude Oliva, and Hanspeter Pfister. 2013. What makes a visualization memorable? *IEEE Transactions on Visualization and Computer Graphics (TVCG)* 19, 12 (2013), 2306–2315.

[13] Matthew Michael Botvinick. 2012. Hierarchical reinforcement learning and decision making. *Current Opinion in Neurobiology* 22, 6 (2012), 956–962.

[14] Dirk Brockmann and Theo Geisel. 2000. The ecology of gaze shifts. *Neurocomputing* 32 (2000), 643–650.

[15] Anthony Brohan, Yevgen Chebotar, Chelsea Finn, Karol Hausman, Alexander Herzog, Daniel Ho, Julian Ibarz, Alex Irpan, Eric Jang, Ryan Julian, et al. 2023. Do as I can, not as I say: Grounding language in robotic affordances. In *Conference on Robot Learning*. PMLR, 287–318.

[16] Victor Cărbune, Hassan Mansoor, Fangyu Liu, Rahul Aralikatte, Gilles Baechler, Jindong Chen, and Abhanshu Sharma. 2024. Chart-based reasoning: Transferring capabilities from LLMs to VLMs. In *Findings of the Association for Computational Linguistics: NAACL 2024*. 989–1004.

[17] Suyog Chandramouli, Danqing Shi, Aini Putkonen, Sebastiaan De Peuter, Shanshan Zhang, Jussi Jokinen, Andrew Howes, and Antti Oulasvirta. 2024. A workflow for building computationally rational models of human behavior. *Computational Brain & Behavior* 7, 3 (2024), 399–419.

[18] Shi Chen, Ming Jiang, Jinhui Yang, and Qi Zhao. 2020. AiR: Attention with Reasoning capability. In *Computer Vision – ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part I 16*. Springer, 91–107.

[19] Xianyu Chen, Ming Jiang, and Qi Zhao. 2021. Predicting human scanpaths in visual question answering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 10876–10885.

[20] Antoine Coutrot, Janet H Hsiao, and Antoni B Chan. 2018. Scanpath modeling and classification with hidden Markov models. *Behavior Research Methods (BRM)* 50, 1 (2018), 362–379.

[21] Murtaza Dalal, Tarun Chiruvolu, Devendra Chaplot, and Ruslan Salakhutdinov. 2024. Plan-Seq-Learn: Language model guided RL for solving long horizon robotics tasks. *International Conference on Learning Representations*.

[22] Gautier Drusch, JC Bastien, and Stéfane Paris. 2014. Analysing eye-tracking data: From scanpaths and heatmaps to the dynamic visualisation of areas of interest. In *Advances in Science, Technology, Higher Education and Society in the Conceptual Age: STHESCA*, Marek Tadeusz (Ed.).

[23] Andrew T Duchowski. 2018. Gaze-based interaction: A 30 year retrospective. *Computers & Graphics* 73 (2018), 59–69.

[24] Manfred Eppe, Christian Gumbsch, Matthias Kerzel, Phuong DH Nguyen, Martin V Butz, and Stefan Wermter. 2022. Intelligent problem-solving as integrated hierarchical reinforcement learning. *Nature Machine Intelligence* 4, 1 (2022), 11–20.

[25] Camilo Fosco, Vincent Casser, Amish Kumar Bedi, Peter O'Donovan, Aaron Hertzmann, and Zoya Bylinskii. 2020. Predicting visual importance across graphic design types. In *Proceedings of the 33rd Annual ACM Symposium on User Interface Software and Technology (UIST)*. 249–260.

[26] Michael J Frank and David Badre. 2012. Mechanisms of hierarchical reinforcement learning in corticostriatal circuits 1: Computational analysis. *Cerebral Cortex* 22, 3 (2012), 509–526.

[27] Joseph Goldberg and Jonathan Helfman. 2011. Eye tracking for visualization evaluation: Reading values on linear versus radial graphs. *Information Visualization* 10, 3 (2011), 182–195.

[28] Joseph H Goldberg and Jonathan I Helfman. 2010. Comparing information graphics: A critical look at eye tracking. In *Proceedings of the 3rd BELIV'10 Workshop: BEyond time and errors: novel evaLuation methods for Information Visualization*. 71–78.

[29] Yi Guo, Nan Cao, Ligan Cai, Yanqiu Wu, Daniel Weiskopf, Danqing Shi, and Qing Chen. 2023. Datamator: An authoring tool for creating datamations via data query decomposition. *Applied Sciences* 13, 17 (2023), 9709.

[30] Yi Guo, Danqing Shi, Mingjuan Guo, Yanqiu Wu, Nan Cao, and Qing Chen. 2024. Talk2data: A natural language interface for exploratory visual analysis via question decomposition. *ACM Transactions on Interactive Intelligent Systems* 14, 2 (2024), 1–24.

[31] Christopher Healey and James Enns. 2011. Attention and visual memory in visualization and computer graphics. *IEEE Transactions on Visualization and Computer Graphics* 18, 7 (2011), 1170–1188.

[32] Stacie L Hibino. 1999. Task analysis for information visualization. In *International Conference on Advances in Visual Information Systems*. Springer, 139–146.

[33] Weidong Huang. 2007. Using eye tracking to investigate graph layout effects. In *2007 6th International Asia-Pacific Symposium on Visualization*. IEEE, 97–100.

[34] Wenlong Huang, Pieter Abbeel, Deepak Pathak, and Igor Mordatch. 2022. Language models as zero-shot planners: Extracting actionable knowledge for embodied agents. In *International Conference on Machine Learning*. PMLR, 9118–9147.

[35] Wenlong Huang, Fei Xia, Ted Xiao, Harris Chan, Jacky Liang, Pete Florence, Andy Zeng, Jonathan Tompson, Igor Mordatch, Yevgen Chebotar, et al. 2023. Inner monologue: Embodied reasoning through planning with language models. In *Conference on Robot Learning*. PMLR, 1769–1782.

[36] Laurent Itti and Christof Koch. 2000. A saliency-based search mechanism for overt and covert shifts of visual attention. *Vision Research* 40, 10–12 (2000), 1489–1506.

[37] Laurent Itti, Christof Koch, and Ernst Niebur. 1998. A model of saliency-based visual attention for rapid scene analysis. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)* 20, 11 (1998), 1254–1259.

[38] Yue Jiang, Zixin Guo, Hamed Rezazadegan Tavakoli, Luis A Leiva, and Antti Oulasvirta. 2024. EyeFormer: Predicting personalized scanpaths with Transformer-guided reinforcement learning. In *UIST '24: Proceedings of the 37th Annual ACM Symposium on User Interface Software and Technology* (Pittsburgh, PA, USA). ACM, New York, NY, USA, Article 47. https://doi.org/10.1145/3654777.3676436

[39] Jussi PP Jokinen, Zhenxin Wang, Sayan Sarcar, Antti Oulasvirta, and Xiangshi Ren. 2020. Adaptive feature guidance: Modelling visual search with graphical layouts. *International Journal of Human–Computer Studies* 136, Article 102376 (2020).

[40] Florian Kadner, Yannik Keller, and Constantin Rothkopf. 2021. AdaptiFont: Increasing individuals' reading speed with a generative font model and Bayesian optimization. In *CHI '21: Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*. Article 585.

[41] Matthias Kümmerer, Matthias Bethge, and Thomas SA Wallis. 2022. DeepGaze III: Modeling free-viewing human scanpaths with deep learning. *Journal of Vision* 22, 5, Article 7 (2022).

[42] Sébastien Lallé, Tiffany Wu, and Cristina Conati. 2020. Gaze-driven links for magazine style narrative visualizations. In *2020 IEEE Visualization Conference (VIS)*. 166–170.

[43] Boyi Li, Philipp Wu, Pieter Abbeel, and Jitendra Malik. 2023. Interactive task planning with language models. *arXiv preprint arXiv:2310.10645* (2023).

[44] Zhi Li, Yu-Jung Ko, Aini Putkonen, Shirin Feiz, Vikas Ashok, IV Ramakrishnan, Antti Oulasvirta, and Xiaojun Bi. 2023. Modeling touch-based menu selection performance of blind users via reinforcement learning. In *CHI '23: Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*. Article 357.

[45] Jacky Liang, Wenlong Huang, Fei Xia, Peng Xu, Karol Hausman, Brian Ichter, Pete Florence, and Andy Zeng. 2023. Code as policies: Language model programs for embodied control. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 9493–9500.

[46] Zhicheng Liu and John Stasko. 2010. Mental models, visual reasoning and interaction in information visualization: A top-down perspective. *IEEE Transactions on Visualization and Computer Graphics* 16, 6 (2010), 999–1008.

[47] Geoffrey R Loftus and Elizabeth F Loftus. 2019. *Human memory: The processing of information*. Psychology Press.

[48] Ahmed Masry, Xuan Long Do, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. 2022. ChartQA: A benchmark for question answering about charts with visual and logical reasoning. In *Findings of the Association for Computational Linguistics: ACL 2022*. 2263–2279.

[49] Laura E Matzen, Michael J Haass, Kristin M Divis, Zhiyuan Wang, and Andrew T Wilson. 2017. Data visualization saliency model: A tool for evaluating abstract data visualizations. *IEEE Transactions on Visualization and Computer Graphics (TVCG)* 24, 1 (2017), 563–573.

[50] Sounak Mondal, Zhibo Yang, Seoyoung Ahn, Dimitris Samaras, Gregory Zelinsky, and Minh Hoai. 2023. Gazeformer: Scalable, effective and fast prediction of goal-directed human attention. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 1441–1450.

[51] Roderick Murray-Smith, Antti Oulasvirta, Andrew Howes, Jörg Müller, Aleksi Ikkala, Miroslav Bachinski, Arthur Fleig, Florian Fischer, and Markus Klar. 2022. What simulation can do for HCI research. *Interactions* 29, 6 (2022), 48–53.

[52] Saul B Needleman and Christian D Wunsch. 1970. A general method applicable to the search for similarities in the amino acid sequence of two proteins. *Journal of Molecular Biology* 48, 3 (1970), 443–453.

[53] Truong-Huy D Nguyen, Magy Seif El-Nasr, and Derek M Isaacowitz. 2015. Interactive visualization for understanding of attention patterns. In *Workshop on Eye Tracking and Visualization (ETVIS)*. 23–39.

[54] Antti Oulasvirta, Jussi PP Jokinen, and Andrew Howes. 2022. Computational rationality as a theory of interaction. In *CHI '22: Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems*. Article 359.

[55] Richard W Pew. 1966. Acquisition of hierarchical control over the temporal organization of a skill. *Journal of Experimental Psychology* 71, 5 (1966), 764–771.

[56] Patrik Polatsek, Manuela Waldner, Ivan Viola, Peter Kapec, and Wanda Benesova. 2018. Exploring visual attention and saliency modeling for task-based visual analysis. *Computers & Graphics* 72 (2018), 26–38.

[57] Maryam Rezaie, Melanie Tory, and Sheelagh Carpendale. 2024. Struggles and strategies in understanding information visualizations. *IEEE Transactions on Visualization and Computer Graphics* (2024).

[58] Arvind Satyanarayan, Dominik Moritz, Kanit Wongsuphasawat, and Jeffrey Heer. 2016. Vega-Lite: A grammar of interactive graphics. *IEEE Transactions on Visualization and Computer Graphics* 23, 1 (2016), 341–350.

[59] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347* (2017).

[60] Hans-Jörg Schulz, Thomas Nocke, Magnus Heitzler, and Heidrun Schumann. 2013. A design space of visualization tasks. *IEEE Transactions on Visualization and Computer Graphics* 19, 12 (2013), 2366–2375.

[61] Danqing Shi, Antti Oulasvirta, Tino Weinkauf, and Nan Cao. 2024. Understanding and automating graphical annotations on animated scatterplots. In *2024 IEEE 17th Pacific Visualization Conference (PacificVis)*. IEEE, 212–221.

[62] Danqing Shi, Yang Shi, Xinyue Xu, Nan Chen, Siwei Fu, Hongjin Wu, and Nan Cao. 2019. Task-oriented optimal sequencing of visualization charts. In *2019 IEEE Visualization in Data Science (VDS)*. IEEE, 58–66.

[63] Danqing Shi, Fuling Sun, Xinyue Xu, Xingyu Lan, David Gotz, and Nan Cao. 2021. AutoClips: An automatic approach to video generation from data facts. 40, 3 (2021), 495–505.

[64] Danqing Shi, Xinyue Xu, Fuling Sun, Yang Shi, and Nan Cao. 2020. Calliope: Automatic visual data story generation from a spreadsheet. *IEEE Transactions on Visualization and Computer Graphics* 27, 2 (2020), 453–463.

[65] Danqing Shi, Yujun Zhu, Jussi PP Jokinen, Aditya Acharya, Aini Putkonen, Shumin Zhai, and Antti Oulasvirta. 2024. CRTypist: Simulating touchscreen typing behavior via computational rationality. In *CHI '24: Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*. Article 942.

[66] Sungbok Shin, Sunghyo Chung, Sanghyun Hong, and Niklas Elmqvist. 2022. A scanner deeply: Predicting gaze heatmaps on visualizations using crowdsourced eye movement data. *IEEE Transactions on Visualization and Computer Graphics (TVCG)* 29, 1 (2022), 396–406.

[67] Sungbok Shin, Sanghyun Hong, and Niklas Elmqvist. 2023. Perceptual Pat: A virtual human visual system for iterative visualization design. In *CHI '23: Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*. Article 811.

[68] Raghuraj Singh, Chandra Shekhar Yadav, Prabhat Verma, and Vibhash Yadav. 2010. Optical character recognition (OCR) for printed Devnagari script using artificial neural network. *International Journal of Computer Science & Communication* 1, 1 (2010), 91–95.

[69] Ekta Sood, Fabian Kögel, Philipp Müller, Dominike Thomas, Mihai Bâce, and Andreas Bulling. 2023. Multimodal integration of human-like attention in visual question answering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2648–2658.

[70] Arjun Srinivasan, Steven M Drucker, Alex Endert, and John Stasko. 2018. Augmenting visualizations with interactive data facts to facilitate interpretation and communication. *IEEE Transactions on Visualization and Computer Graphics* 25, 1 (2018), 672–681.

[71] Xiangjie Sui, Yuming Fang, Hanwei Zhu, Shiqi Wang, and Zhou Wang. 2023. ScanDMM: A deep Markov model of scanpath prediction for 360deg images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition(CVPR)*. 6989–6999.

[72] Wanjie Sun, Zhenzhong Chen, and Feng Wu. 2019. Visual scanpath prediction using IOR-ROI recurrent mixture density network. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 43, 6 (2019), 2101–2118.

[73] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. LLaMA: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971* (2023).

[74] Sai H Vemprala, Rogerio Bonatti, Arthur Bucker, and Ashish Kapoor. 2024. ChatGPT for robotics: Design principles and model abilities. *IEEE Access* 12 (2024).

[75] Ashish Verma and Debashis Sen. 2019. HMM-based convolutional LSTM for visual scanpath prediction. In *2019 27th European Signal Processing Conference (EUSIPCO)*. 1–5.

[76] Yao Wang, Mihai Bâce, and Andreas Bulling. 2024. Scanpath prediction on information visualisations. *IEEE Transactions on Visualization and Computer Graphics (TVCG)* 30, 7 (2024), 3902–3914.

[77] Yao Wang, Yue Jiang, Zhiming Hu, Constantin Ruhdorfer, Mihai Bâce, and Andreas Bulling. 2024. VisRecall++: Analysing and predicting visualisation recallability from gaze behaviour. *Proceedings of the ACM on Human–Computer Interaction* 8, ETRA, Article 339 (2024).

[78] Yao Wang, Maurice Koch, Mihai Bâce, Daniel Weiskopf, and Andreas Bulling. 2022. Impact of gaze uncertainty on AOIs in information visualisations. In *Proceedings of the ACM International Symposium on Eye Tracking Research and Applications (ETRA)*. Article 60.

[79] Yixiu Wang, Bin Wang, Xiaofeng Wu, and Liming Zhang. 2017. Scanpath estimation based on foveated image saliency. *Cognitive Processing* 18, 1 (2017), 87–95.

[80] Yao Wang, Weitian Wang, Abdullah Abdelhafez, Mayar Elfares, Zhiming Hu, Mihai Bâce, and Andreas Bulling. 2024. SalChartQA: Question-driven saliency on information visualisations. In *CHI '24: Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*. Article 763.

[81] Hongjin Wu, Danqing Shi, Nan Chen, Yang Shi, Zhuochen Jin, and Nan Cao. 2020. VisAct: a visualization design system based on semantic actions. *Journal of Visualization* 23 (2020), 339–352.

[82] Zhibo Yang, Lihan Huang, Yupei Chen, Zijun Wei, Seoyoung Ahn, Gregory Zelinsky, Dimitris Samaras, and Minh Hoai. 2020. Predicting goal-directed human attention using inverse reinforcement learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 193–202.

[83] Zhibo Yang, Sounak Mondal, Seoyoung Ahn, Ruoyu Xue, Gregory Zelinsky, Minh Hoai, and Dimitris Samaras. 2024. Unifying top-down and bottom-up scanpath prediction using Transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 1683–1693.