# VisQA: Quantifying Information Visualisation Recallability via Question Answering

Yao Wang, Chuhan Jiao, Mihai Bâce, Andreas Bulling

**Abstract**—Despite its importance for assessing the effectiveness of communicating information visually, fine-grained recallability of information visualisations has not been studied quantitatively so far. In this work we propose a visual question answering (VQA) paradigm to study visualisation recallability and present VisQA — a novel VQA dataset consisting of 200 visualisations that are annotated with crowd-sourced human (N = 305) recallability scores obtained from 1,000 questions from five question types. Furthermore, we present the first computational method to predict recallability of different visualisation elements, such as the title or specific data values. We report detailed analyses of our method on VisQA and demonstrate that it outperforms several baselines in overall recallability and FE-, F-, RV-, and U-question recallability. We further demonstrate one possible application of our method: recommending the visualisation type that maximises user recallability for a given data source. Taken together, our work makes fundamental contributions towards a new generation of methods to assist designers in optimising visualisations.

**Index Terms**—Information Visualisation, Recallability, Visual Question Answering, Memorability, Machine Learning

✦

## 1 INTRODUCTION

Memorability is an intrinsic, global, and stimulus-driven perceptual property that is important for better comprehension of visual stimuli [1, 2]. A growing body of work has studied image recognisability – one of the most fundamental attributes of memorability, both from a perceptual [1, 3] and a computational [4, 5] perspective. Recognisability has also been studied on information visualisations and previous work has revealed specific attributes that make visualisations memorable [6]. Recognisability measures whether a visualisation looks familiar or novel [3]. A visualisation that has unique features may stand out more and may therefore be more memorable. However, recognisability does not capture how effective a visualisation is in conveying information to observers. Other works have therefore studied *recallability* – a concept that goes beyond memorability, yet is complementary to it [7], by quantifying *what* viewers remember from a visualisation [8]. Despite its importance and potential for designing better information visualisations, a deeper understanding of which characteristics of visualisations influence recallability, and in which way, is currently missing.

Current methods to assess recallability rely on visualisation experts to assign a qualitative score to self-reported free-text descriptions of viewers [7]. This approach is cumbersome and only provides a single score representing overall recallability while hiding the contribution of individual visualisation characteristics. While Borkin et al. [7] noted the importance of titles for recallability on visualisations, Polatsek et al. [9] conducted three low-level analytical tasks, focusing on visual elements with extrema, or specific values. These works inspired us to quantify visualisations'

recallability by looking into specific types of visualisation elements, such as the title, elements with extrema, or distinct data points.

To quantify recallability, we propose to adopt a question-answering paradigm, similar to visual question answering (VQA) [10] that has become widely popular in computer vision. VQA involves computational models in reasoning and correctly answering questions about images. While originally introduced for natural images [10, 11], VQA was also explored for information visualisations [12]. One follow-up work collected human performance values for the DVQA dataset by crowd-sourcing [13]. Inspired by this, we evaluate the performance of observers in answering questions about images correctly and use their performance as a subjective measure of information visualisation recallability.

In this work, to quantify fine-grained recallability of information visualisations, we design and execute a VQA-based study to collect VisQA: a novel visualisation dataset with 200 visualisations, which contains 1000 high quality questions annotated by visualisation experts and crowd-sourced human recallability scores. Our work is inspired by and extends prior task taxonomy on visualisations [9] to define fine-grained recallability scores [14] through five question types: identifying the title or theme, finding extrema, filtering data elements, retrieving values, and understanding structure (subsection 3.1). Through our analyses of VisQA, we make several interesting findings: the highest recallability across question types occurs in questions that are about the title or the general theme (T-question), which is significantly higher than other question types. Our replication study of recognition accuracy aligns well with the previous memorability studies [6, 7], and we conclude that there are no such visualisations with high recallability and low memorability. This finding suggests that recallability is more descriptive than memorability and more challenging to predict on information visualisations. Based on VisQA, we further present RecallNet, a novel method based on

---

- *Yao Wang, Mihai Bâce, and Andreas Bulling are with the Institute for Visualisation and Interactive Systems, University of Stuttgart, Germany, E-mail: {yao.wang, mihai.bace, andreas.bulling}@vis.uni-stuttgart.de.*
- *Chuhan Jiao is with Aalto University, E-mail: chuhan.jiao@aalto.fi*
- *Yao Wang is the corresponding author.*

convolutional neural networks (CNNs) to predict one overall and five fine-grained recallability scores, one for each question type. Finally, we prototype a novel application for visualisation type recommendations that maximises user recallability. Triplets of information visualisations are created with minimised content changes across visualisation types. The recallability scores on the visualisation triplets are then collected. Through a user study, we demonstrate that the prediction from our RecallNet not only maximises user recallability but also agrees with the preferences from scientific researchers in three out of four visualisation triplets.

Our contribution is threefold: (1) We adapt a visual question answering (VQA) paradigm to quantify fine-grained recallability of information visualisations. (2) We collect VisQA, a novel visualisation dataset with human recallability scores (N = 305) from 1000 questions and five question types. (3) We propose a computational model that predicts fine-grained recallability of visualisations and demonstrates how our model can be used to automatically recommend a visualisation type that increases recallability. As such, our work points the way towards new methods and tools to create more effective information visualisations.

## 2 RELATED WORK

Our work is related to previous works on 1) image memorability, 2) perception and memorability of visualisations, and 3) visualisation visual question answering (VQA) datasets.

### 2.1 Image Memorability

A pioneering study [3] reported a strong capability of humans to recognise what they have seen before even up to 10,000 images, which is denoted as "image recognition memory". The following studies have demonstrated that memorability is an observer-independent property, which only depends on images [15, 16]. Furthermore, previous studies have proven that memorability could be reliably quantified for individual images by asking subjects to report whether images are novel or familiar [4, 17]. Large-scale memorability datasets have been collected for natural images, such as SUN-Mem [4], Figrim [18] and LaMem [5]. With the rise of deep learning, deep convolutional neural networks were proposed as computational methods to predict image memorability [5, 19, 20]. Recent work also integrated visual attention into the memorability prediction model [21]. Meanwhile, recallability is a complementary memory task to visual recognition [22], which requires subjects to view images and then recall what they have seen [23]. One previous work found out that sketch-based methodologies can improve the recall of a sampling distribution from an experiment [24]. Several recent studies are consistent with the conclusion that image memorability variation for recognition and recall tasks may be distinct [8, 25]. Based on this, our work is the first to better understand the recallability characteristics and the factors that influence it on information visualisations.

### 2.2 Perception and Memorability of Visualisations

Pioneering works in the visualisation community have examined how different data types and tasks influence human perception [26, 27, 28]. Bateman et al. claimed that the overembellishment (i.e., "chart junk") improves recognisability but is not essential for understanding the visualisation [29]. This triggered a series of studies evaluating the impact of style on memorability and comprehensibility [30, 31, 32]. The effect of specific factors or components on recall memory has been investigated, such as interaction [33], prior knowledge [34], title [7, 35] and text redundancy [7]. Borkin et al. [6] studied visualisation memorability on the MASSVIS dataset, and their follow-up work [7] further conducted online crowd-sourcing studies to quantify both recognisability and recallability. There are two main drawbacks to the previous recallability quantification procedure. Firstly, the method used to recall quality annotations is subjective and cumbersome. In addition, visualisation experts are necessary to attribute these scores. Secondly, the quality score scale with only four possible values is too coarse to represent a visualisation. To overcome these limitations, we introduce visual question answering (VQA) as a powerful paradigm to quantify the recallability of information visualisations. Through multiple questions and answers on different visualisation characteristics, we propose a novel computational model to predict not only overall but also fine-grained recallability based on five different question types.

### 2.3 Visualisation Visual Question Answering (VQA) Dataset

The visual question answering (VQA) task [10] proposed in the field of computer vision has triggered many follow-up studies and applications [11, 36]. Despite the importance of information visualisations, the visualisation VQA datasets have only been proposed in recent years. FigureQA [12] was the first visualisation VQA dataset. Images were plotted in simple and fully synthesised visualisations in five visualisation classes, along with polar questions. DVQA [37] is a dataset focused specifically on the problem of visual reasoning on bar charts, which is used as a corpus for the topic of chart QA. PlotQA [38] and LEAF-QA [39] synthesised their question-answer pairs based on crowd-sourced question templates from real-world data sources to increase the variety. Kafle et al. [13] collected human performance values for the DVQA dataset using crowd-sourcing. As a conclusion, VQA has not been used for memorability studies yet, and current visualisation VQA datasets are synthesised from simple templates with limited content, making it a distance away from real world visualisations. However, VQA provides an interesting means to get fine-grained annotations and insights into recallability. In our work, we evaluate and obtain recallability scores by asking users questions and validating their answers. Therefore, we present the design of our novel adaptation of a VQA-based study on information visualisations and our novel VisQA dataset in the next section.

## 3 VISQA DATASET

The currently available recallability scores on the visualisation dataset MASSVIS [6, 7] are annotated from free-text
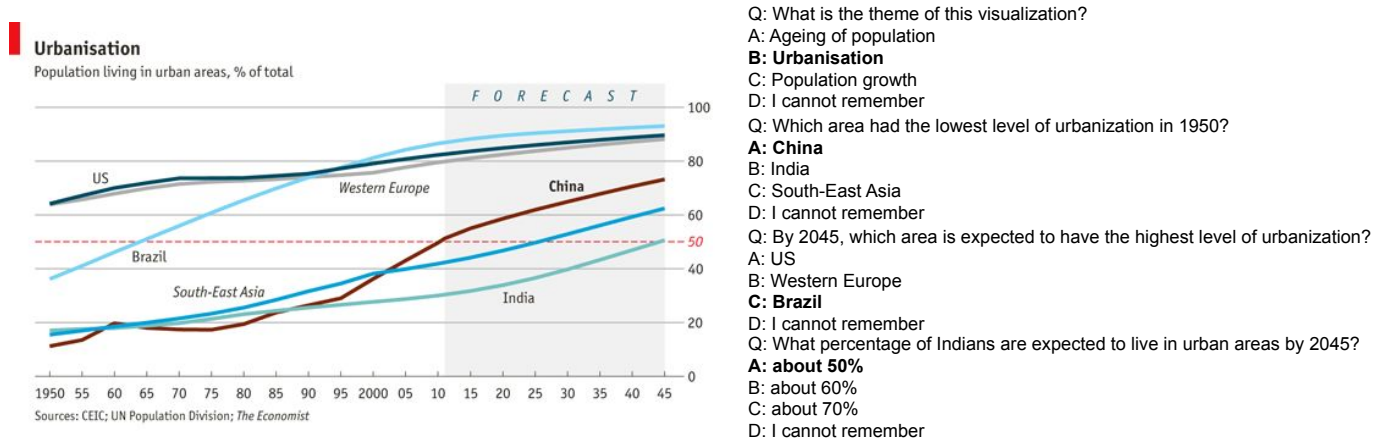
Q: What is the theme of this visualization?
A: Ageing of population
**B: Urbanisation**
C: Population growth
D: I cannot remember
Q: Which area had the lowest level of urbanization in 1950?
**A: China**
B: India
C: South-East Asia
D: I cannot remember
Q: By 2045, which area is expected to have the highest level of urbanization?
A: US
B: Western Europe
**C: Brazil**
D: I cannot remember
Q: What percentage of Indians are expected to live in urban areas by 2045?
**A: about 50%**
B: about 60%
C: about 70%
D: I cannot remember

Fig. 1: Sample visualisation with multiple-choice questions from VisQA. Five types of questions were designed by experts, which are questions regarding the title (T-questions), understanding structure or trend (U-questions), finding extrema (FE-questions), filtering elements (F-questions) and retrieving values (RV-questions). Each figure has at least two question types. Image sourced from MASSVIS [6].

descriptions. However, its procedure to quantify recallability is coarse and cumbersome. Meanwhile, Visual Question-Answering (VQA) datasets [10] selectively target elements of visualisations in different question-answer pairs, making it a suitable setting to quantify memorability objectively and efficiently. Under the VQA paradigm, different tasks can be represented as different types of questions to viewers, and consequently, recallability is quantified by the accuracy in answering those questions.

Towards quantifying recallability, we propose the Visualisation Recallability Question Answering Dataset (VisQA) — a VQA dataset with 200 real-world information visualisations and contains crowd-sourced human recallability scores (N = 305) obtained from 1,000 questions in five question types (see Figure 1). Visualisations in our dataset are mainly sourced from the MASSVIS dataset [6] to enable better alignment with prior works on this topic. The recognisability scores are also collected to replicate the previous memorability studies [6, 7].

### 3.1 Visualisation Collection and Question Types

We randomly selected a subset of 200 visualisations from the MASSVIS dataset [6]. Notably, we excluded all info-graphics in our collection, since infographics have the highest recognisability and recallability compared to all other types of visualisations [7]. However, scatter plots represent only 5 % of the sampled subset. Therefore, we collected 20 additional scatter plot visualisation by crawling the web through search engines (Google, Bing) using the keyword of *scatter plots*. Then, we replaced some bar plots with the web-crawled plots to balance the visualisation type classes. The final distribution of visualisation types is: 60 bar plots, 45 line plots, 28 scatter plots, 19 pie plots, 19 tables and 19 others. Those visualisations that don't belong to any of the first five types are categorised as *others*, including box charts, isotype charts, or other complex visualisations.

VisQA contains five types of questions: T-questions, U-questions, FE-questions, F-questions, and RV-questions. T-questions are questions regarding the title or the visualisation theme, and U-questions are about understanding the

plot structure [38] or the general trend [39]. The remaining three question types correspond to three low-level analytical tasks introduced in [9], which are finding an extremum attribute value (FE-questions), filtering data elements based on specific criteria (F-questions) and retrieving values for a specific data element (RV-questions).

All question answering data were created by five data visualisation experts. They were asked to provide five questions per visualisation, and every visualisation has at least two question types. Each question corresponds to four possible answer options. Only one option is correct, two other options are choices with similar, yet incorrect answers, and the last option is always "I cannot remember". See supplementary material for question examples. All annotations were saved separately in standard JSON files for each visualisation. There are 193, 150, 178, 99, 64 visualisations in VisQA that have at least one T-, FE-, F-, RV-, and U-question, respectively.

**T-question.** T-questions are about the title or the general theme of the plot and do not require any reasoning. Example questions: *What is the title of the visualisation?*, *What is the theme of the visualisation?* For the incorrect choices in T-questions, we either replaced keywords or phrases with words of similar, but different meanings, such as changing *car thefts* to *car accidents* or *car manufacturers*, or used titles from other visualisations, such as using *Covered Transactions by Sector and Year, 2009-2011* and *HIV Prevalence in Women Aged 15-49 Years by Region, 1990-2007* as incorrect choices for *Covered Transactions by Sector and Year, 2009-2011*.

**FE-question.** These are questions about finding extreme values in the visualisation that fulfil certain conditions, without asking any exact numbers. Example questions: *Which area had the lowest level of urbanization in 1950?*, and *Which particle is the latest discovered?* For FE-questions, we used other elements that appeared in the visualisation as incorrect answer choices. For example, *India* and *South-East Asia* were the incorrect alternative choices for *China* in the question *Which area had the lowest level of urbanization in 1950?* (see Figure 1).

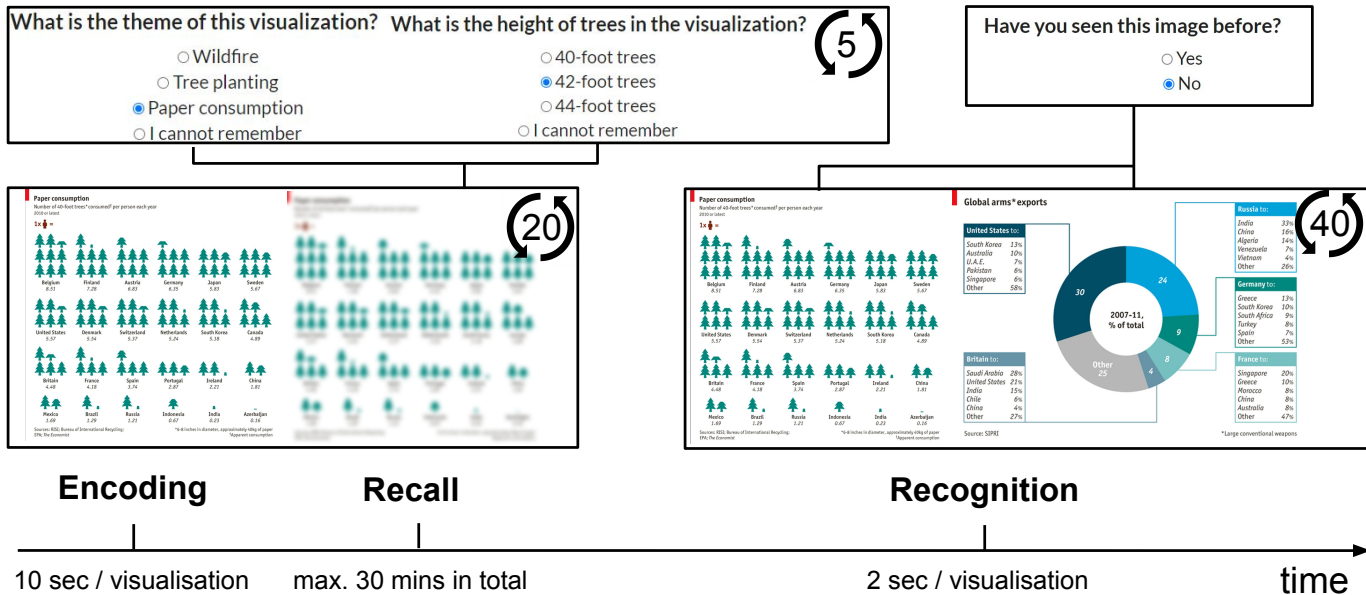**F-question.** These are questions about filtering data ele-

Fig. 2: Experiment design. From Left to Right: Visualisations are shown to viewers for 10 seconds in the "encoding" phase. In the "recall" phase, visualisations are blurred and each have a multiple-choice question next to it with a single correct answer. Finally, visualisations are shown to viewers for 2 seconds in the "recognition" phase.

ments based on specific criteria. Example questions: *Which particle is Bosons?* and *What is the source of the data?* For F-questions, we either changed keywords to their synonyms, or used other elements that appeared in visualisations as incorrect alternative choices, such as using *Electron* and *Muon* for *Photon* in the question *Which particle is Bosons?*.

**RV-question.** These are questions about retrieving a specific value located in the plot. Those RV-questions in combination with FE-questions are all categorised in this type. Example questions: *What is the maximum percentage of aid allocated?* and *What percentage of Indians are expected to live in urban areas by 2045?* (see Figure 1). Example incorrect choices: *about 60 %* and *about 70 %* for *What percentage of Indians are expected to live in urban areas by 2045?*, and the correct answer is *about 50 %*.

**U-question.** These are questions about understanding the structure or the trend of a visualisation. Example questions: *What does the purple curve represent?* and *What decreases as time goes by?* Example incorrect choices: for structure questions, other elements appearing in the visualisation are used, such as using *Red* and *Blue* as incorrect choices for *Green* in the question *What color stands for Residents?* As for questions about understanding trends, the choices are *increasing*, *decreasing* and *almost the same*.

### 3.2 Crowd-sourcing Study Set-up & Participants

Our study design is illustrated in Figure 2. In the encoding phase of our study, study participants were shown a sequence of visualisations for 10 seconds each, which is similar to the prior memorability study [6]. We asked participants to memorise as much of the information presented in each visualisation as possible. In the recall phase, we showed participants the blurred image of the first visualisation with a single choice question. The following question would be shown only if they clicked the next button, and

they could not return to the previous question. This setting was to avoid providing hints in upcoming questions. After answering all five questions, the process for the second visualisation was the same. Then, the recognition phase involved an online memorability game similar to the prior work of Borkin et al. [6]. Study participants were presented with a sequence of images, and they had to select if they had seen this visualisation before. In each Human Intelligence Task (HIT), 40 blurred images were shown for 2 seconds each. The images in the recognition phase contained 20 visualisations that were the same in the recall phase, and 20 fillers from a different group. Finally, participants were asked to provide anonymous feedback on the study design in a questionnaire.

To support the VQA setting, we implemented question answering procedures in a web application. We blurred all visualisations with Gaussian filters, kernel_size adaptive to the image resolution, ranging from 5 to 24. We then integrated our VQA application into an existing crowd-sourcing toolbox that worked well with Amazon Mechanical Turk platform [40]. We deployed our experiment on Amazon's Mechanical Turk (MTurk) platform to collect recallability and recognisability scores on all 200 visualisations, splitting them randomly into ten groups of 20 visualisations per HIT. MTurk workers could participate in multiple groups. To participate in one of our HITs, a worker had to be a Master Worker approved by MTurk as a quality check. Master Workers are top workers rated by the MTurk who have consistently demonstrated high quality works. Workers were paid \$ 4.00 for completing each HIT. To ensure data quality, we filtered out 467 HITs (N = 305 workers) if the answers were all "Yes" or "No" in the recognition task. For each visualisation, we received an average of 40.4 (SD = 16.9) valid responses. The 305 workers were distributed in various educational levels: 8.2 % two-year degree, 56.9 % four-
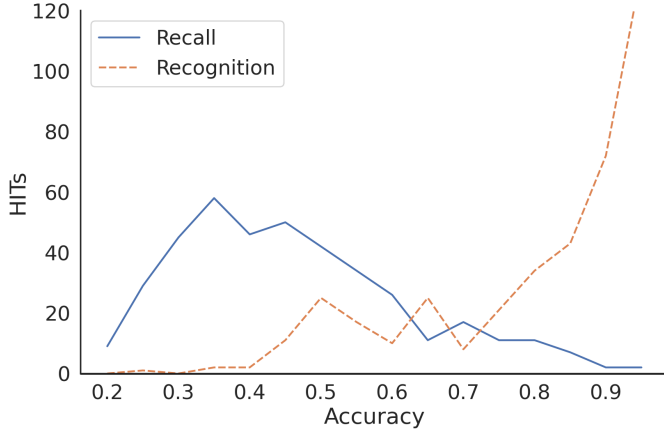
Fig. 3: Recallability and Recognition accuracy over all 404 HITs. Participants can recognise most of the visualisations easily, but they can only answer around half of the questions correctly.
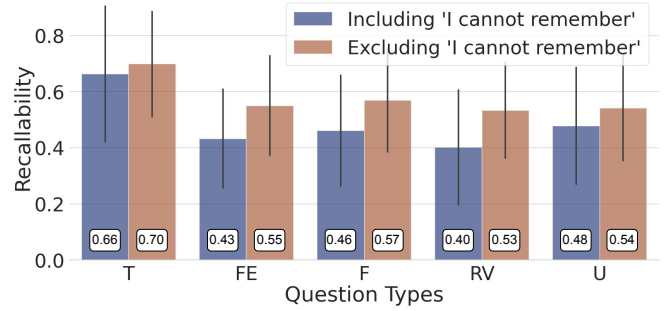


Fig. 4: Recallability scores by question type. T-questions have a significantly higher accuracy compared with all other question types (FE, F, RV and U). 24.7 % of the viewers selected *I cannot remember* in RV-questions, and only 5.1 % of the viewers selected *I cannot remember* in T-questions.

year degree, 22.3 % master's degree or higher, and 12.6 % other/unreported. The age groups were 44.1 % in 25-34, 28.5 % in 35-44, 12.4 % in 45-55 and 9.9 % over 55. In the anonymous feedback form at the end of our study, most workers responded positively and two examples being: "Great self test for capable of memory power" and "After taking survey, I'm really getting interested in learning data plots and visualisations".

### 3.3 Data Analysis

**Recallability formulation.** For each question, we measured the recall accuracy as follows: $Acc = \frac{RA}{RA+WA}$, where $RA$ is the number of correct answers, $WA$ is the number of wrong answers, including the number of *I cannot remember* answers. If we focus on viewers who have selected choices excluding *I cannot remember*, the accuracy can be computed as: $Acc' = \frac{RA}{RA+WA-CNR}$, where CNR stands for the number of *I cannot remember*. Averaging all questions of type $t$ in a visualisation gives us the recallability by question type and is computed as: $Rec_t = \frac{1}{n}\sum_{i=1}^{n} Acc(i), question_i \in t$. By averaging all questions in a visualisation, we have the overall recallability of a visualisation as: $Rec = \frac{1}{n}\sum_{i=1}^{n} Acc(i)$.

**HIT-wise Recallability.** HIT-wise recallability as well as recognition accuracy across HITs (N = 404) are shown in Figure 3. 63.9 % of HITs had a recognition accuracy higher than 0.85, and 34.83 % were higher than 0.95, which shows that our study participants could easily recognise most of the visualisations (M = 0.83). Meanwhile, they could only answer about half of the questions correctly (M = 0.49, t (404) = 30.05, p < 0.001).

**Fine-grained Recallability by Question Type.** Figure 4 illustrates that T-questions have the highest recall accuracy among all question types (average M = 0.66 including *I cannot remember*, and M = 0.69 excluding *I cannot remember*). The accuracy of T-questions is significantly higher than other question types (t (1969) = 18.87, p < 0.001). 24.7 % of viewers selected *I cannot remember* in RV-questions, and 21.4 %, 18.8 %, 11.7 % for FE-, F- and U-questions, respectively. Only 5.1 % of the study participants selected *I cannot remember* in T-questions. We observed a mean proportion of 19.1 %

(SD = 13.0 %) of study participants who selected *I cannot remember* from all visualisations. The lowest proportion is 3 %, while more than 50 % selected *I cannot remember* in seven specific visualisations.

Figure 5 shows visualisations with the most and fewest *I cannot remember* answers from VisQA. We observe that increased visualisation complexity is a common characteristic for visualisations with the most *I cannot remember* answers. The encoding phase in our study only lasted for ten seconds, which might be too short for some complex visualisations.

**Recognisability: a Comparison to Prior Work.** For a comparison to prior work on recognisability [6, 7], we also calculated the memorability (or recognisability) score on VisQA. According to Borkin et al. [6], the hit rate (HR) and false alarm rate (FAR) were computed as: $HR = \frac{HITS}{HITS+MISSES}$, and $FAR = \frac{FA}{FA+CR}$. Then, the recognisability (memorability) of a visualisation was measured as: $d' = Z(HR) - Z(FAR)$, where Z was the inverse cumulative Gaussian distribution. Figure 6 (Left) shows the distribution of the raw HR scores of all visualisations from the recognition phase. Figure 6 (Right) shows the highest and lowest ranked visualisations across recognisability (memorability) and recallability from our VisQA dataset. Visualisations in each quadrant were ranked highest or lowest 15 % among all visualisations.

## 4 COMPUTATIONAL MODEL FOR PREDICTING FINE-GRAINED RECALLABILITY

Our analyses on VisQA yielded several insights on recallability in information visualisations. There are currently no baseline methods, neither for predicting overall recallability nor for fine-grained recallability. Existing computational models only aimed for predicting memorability, also known as recognisability. We extend and build on state-of-the-art architectures from other computer vision tasks, such as semantic segmentation [41, 42] and image classification [43, 44], and use such methods as the backbone of our architecture.

We designed our Recallability Network (RecallNet) with the specific goal of predicting both overall and fine-grained recallability scores in one single model (see Figure 7 for an
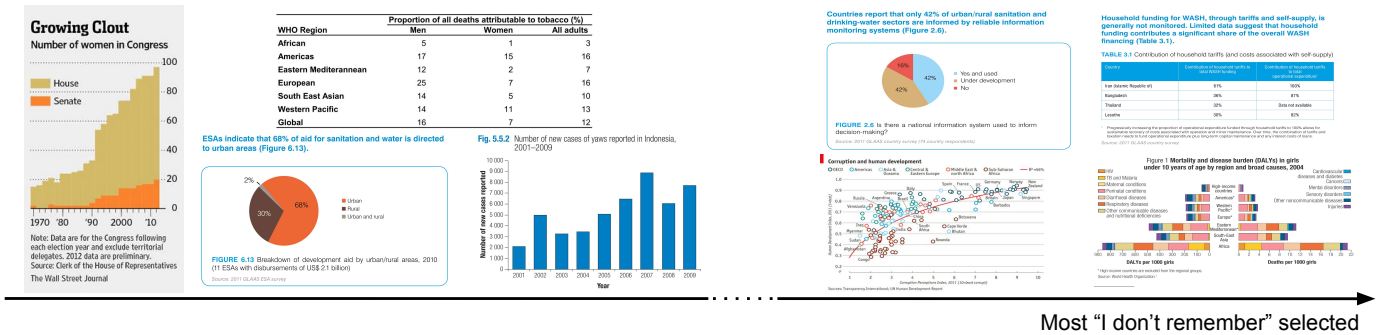
Fig. 5: Example visualisations with the most and fewest answers *I cannot remember* from VisQA. We observed a higher degree of visualisation complexity for those with multiple *I cannot remember* answers.
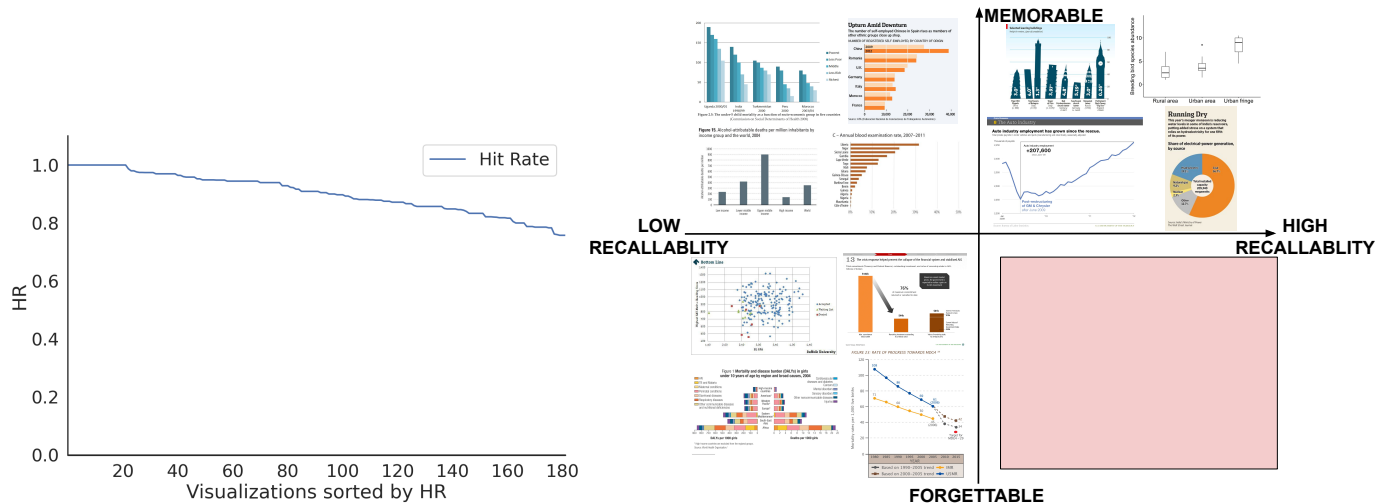


Fig. 6: Left: Raw HR scores of target visualisations from the recognition phase. Right: The highest and lowest ranked visualisations (within 15 %) across recognisability (memorability) and recallability in VisQA. The y-axis represents the memorability score, and the x-axis represents the recallability score computed from overall visualisation question accuracy (independent of question type).

overview). Inspired by UMSI [45], the currently state-of-the-art architecture for visual importance prediction on graphic designs, we employ the Xception [41] model to effectively encode spatial information. Then, a global average pooling layer, a dense layer with 256 neurons, and finally a dense layer with 2 neurons are sequentially connected. One output neuron predicts the general recallability score, and the other one predicts the fine-grained recallability score.

## 5 EXPERIMENTS

### 5.1 Implementation Details & Model Training

We trained RecallNet using weights obtained from the Xception model – which was pretrained on ImageNet [46]. RecallNet was trained with the Adam [47] optimizer with a learning rate of 0.002 and 1:1 Mean Squared Error (MSE) joint loss for the two branches predicting the overall recallability score and the fine-grained recallability score. We averaged all five questions for each image to prepare the ground truth of overall recallability scores. To train our RecallNet to predict fine-grained recallability scores for a certain question type, we only used those visualisations that contained that question type from VisQA. There are 193, 150,

178, 99, and 64 visualisations with at least one T-, FE-, F-, RV-, and U-question, respectively. Five-fold cross-validation was applied to all evaluation processes. All experiments were conducted on a single Nvidia 2060 super GPU with 8GB VRAM.

**Baseline methods.** Since no previous computational models focused on predicting recallability on visualisations, we designed three methods as baselines. We replaced the Xception feature encoder in RecallNet with VGG-16 [43] and ResNet-34 [44] as the first two baselines. The third baseline model is based on UMSI [45], the current state-of-the-art architecture for visual importance prediction. We replaced the decoder in UMSI with a global averaging pooling layer, a dense layer with 256 neurons, and finally a dense layer with 2 neurons — the same final layers as in our model. To be able to use the UMSI model, we annotated the visualisation with their types. We defined six visualisation categories: *scatter plot*, *line plot*, *bar plot*, *pie plot*, *table* and *others*. Visualisations that did not belong to any of the first five types were categorised as *others*. We trained all baseline models for 10 epochs on VisQA starting from ImageNet [46] pretrained weights. We used the Adam [47] optimizer with a learning rate of 0.002 and Mean Squared Error (MSE) loss
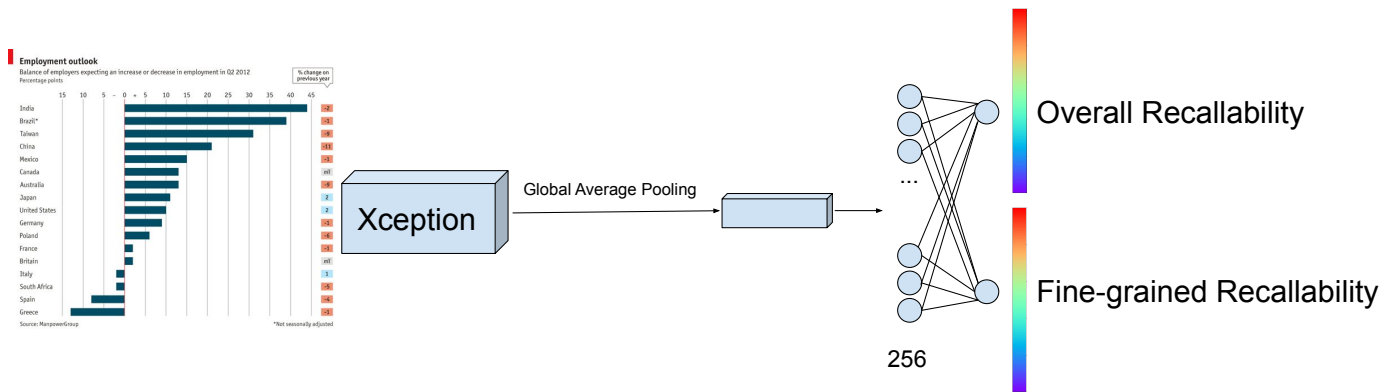
Fig. 7: Method overview. RecallNet leverages the Xception model [41] to effectively encode spatial information. Then, a global average pooling layer, a dense layer with 256 neurons, and finally a dense layer with 2 neurons are sequentially connected. One output neuron predicts the general recallability score, and the other one predicts the fine-grained recallability score.

for training.

### 5.2 Performance Evaluation

We compared the performance of our RecallNet method to the three baselines VGG-16 [43], ResNet-34 [44] and UMSI [45]. Table 1 summarises fine-grained recallability performance on VisQA under 5-fold cross-validation evaluation. The predictions of UMSI [45] collapsed to the same values on every image of validation set, independent of how we adjusted the loss ratio between classification and recallability regression task. This is a typical phenomenon of over-fitting. Our method and the other two baselines do not have over-fitting problems, so we computed the mean squared error for them. Results showed that RecallNet outperformed the baselines under overall recallability and four fine-grained recallability scores, with mean MSE of 0.035 for overall recallability, and 0.021, 0.022, 0.017, 0.043 for FE-, F-, RV-, and U-questions respectively. ResNet-34 was the best performing method for T-questions with a mean MSE of 0.047, while our RecallNet was second with a mean MSE of 0.052.

We also computed the correlation coefficient of recallability scores of RecallNet with those from a previous recallability study [7]. However, we observed non-linear relationships. The correlation coefficient between our overall recallability and their *description quality scores* is -0.013, and between our T-question recallability and their *title mention* is -0.066. The recallability scores from prior work were generated from free-text description without any hints, but our recallability scores were computed from multiple-choice questions with rich context. The non-linear relationships between free-text recallability and multiple-choice recallability suggest that the number of hints provided to viewers are an important factor that may influence recallability.

**Ablation study.** We further carried out an ablation study to investigate how each fine-grained recallability score influences overall recallability (Table 2). In RecallNet, the overall recallability trained with T-questions has the lowest mean squared error of 0.030 and the most stable variance of 0.006. In ResNet-34 [44], the overall recallability trained with RV-questions has the lowest mean squared error of 0.029 and

the most stable variance of 0.008. In VGG-16 [43], the overall recallability trained with T-questions has the lowest mean squared error of 0.037 and the most stable variance of 0.007.

### 5.3 Visualisation Type Recommendation

Visualisation type recommendation is one practical use case for a visualisation recommendation system [48]. Prior research has proposed ways to decide whether line graphs or scatter plots are more suitable for time series data [49]. In our case, we build on our RecallNet to build a prototype implementation that can recommend a visualisation type that maximises recallability. We created four triplets of visualisations from open-source databases[1,2,3]. In each triplet, visualisations have the same data sources but different visualisation types: bar, pie, and line plot, respectively. We used the same colouring scheme, font family, font colour and font size across bar and pie plots to minimise any potential influence of bottom-up saliency [9]. All visualisations in one triplet share the same five multiple-choice questions, and each question belongs to one different type: T, FE, F, RV, and U – introduced in subsection 3.1 (see supplementary materials for figures and questions).

We conducted the same crowd-sourcing study that we used to collect VisQA and recallability scores on the visualisation triplets. All visualisations in four triplets were assigned to three tasks, and every task contained four visualisations, one visualisation from a different triplet. Crowd workers were paid $ 0.80 for completing each HIT. We received 38 (SD = 0.82) valid responses per task. The mean T-question recallability among all triplets was 0.62 (SD = 0.19). The difference when compared to T-question recallability in VisQA was not statistically significant $(t(509) = -1.99, p > 0.95)$. Meanwhile, the mean FE-, F-, RV-, and U-question recallability scores are all significantly higher than VisQA, with t-test scores of 11.03, 8.20, 3.97, 5.66, $p < 0.001$, respectively. Table 3 shows recallability scores on visualisation triplets. We noticed that the recallability of T-questions is

1. https://www.kaggle.com/c/titanic/data
2. https://www.kaggle.com/mohansacharya/graduate-admissions
3. https://data.world/makeovermonday/2021w28

TABLE 1: Fine-grained recallability performance on VisQA under 5-fold cross-validation evaluation. Best results are shown in **bold**, second-best are underlined.

| Methods | Overall ↓ | T ↓ | FE ↓ | F ↓ | RV ↓ | U ↓ |
|---|---|---|---|---|---|---|
| RecallNet (ours) | **0.035 ± 0.005** | 0.052 ± 0.009 | **0.021 ± 0.003** | **0.022 ± 0.004** | **0.017 ± 0.004** | **0.043 ± 0.025** |
| ResNet-34 [44] | 0.043 ± 0.013 | **0.047 ± 0.015** | 0.068 ± 0.024 | 0.070 ± 0.042 | 0.043 ± 0.008 | 0.050 ± 0.018 |
| VGG-16 [43] | 0.036 ± 0.013 | 0.053 ± 0.017 | 0.054 ± 0.019 | 0.076 ± 0.029 | 0.057 ± 0.010 | 0.059 ± 0.025 |
| UMSI [45] | - | - | - | - | - | - |

TABLE 2: Ablation study on the performance of how fine-grained recallability influences the overall recallability. Best results in each row are shown in **bold**.

| Methods | T ↓ | FE ↓ | F ↓ | RV ↓ | U ↓ |
|---|---|---|---|---|---|
| RecallNet (ours) | **0.030 ± 0.006** | 0.079 ± 0.052 | 0.032 ± 0.008 | 0.035 ± 0.013 | 0.172 ± 0.215 |
| ResNet-34 [44] | 0.043 ± 0.013 | 0.078 ± 0.087 | 0.060 ± 0.035 | **0.029 ± 0.008** | 0.033 ± 0.013 |
| VGG-16 [43] | **0.037 ± 0.007** | 0.046 ± 0.022 | 0.041 ± 0.019 | 0.079 ± 0.053 | 0.077 ± 0.011 |

TABLE 3: Fine-grained recallability scores on visualisation triplets. Best results for each visualisation type (pie, bar, or line) are shown in **bold**.

| Type | T ↑ | FE ↑ | F ↑ | RV ↑ | U ↑ | Avg. ↑ |
|---|---|---|---|---|---|---|
| Pie | 0.570 | 0.671 | **0.708** | **0.404** | 0.544 | **0.579** |
| Bar | **0.577** | **0.689** | 0.654 | 0.403 | 0.511 | 0.561 |
| Line | 0.570 | 0.563 | 0.681 | 0.374 | **0.615** | 0.562 |

TABLE 4: Our overall recallability scores and scientific researcher preferences on four pie-bar visualisation pairs. 1 = *pie plot is better* and 5 = *bar plot is better*. The preference is presented as mean and standard deviation. Better results are highlighted in **bold**.

| Type | Overall Recallability ↑ | Preference | Agreed with preference? |
|---|---|---|---|
| Pie | 0.581 | 4.33 (1.25) | ✓ |
| Bar | **0.623** | | |
| Pie | **0.612** | 3.11 (1.73) | ✗ |
| Bar | 0.586 | | |
| Pie | **0.625** | 2.22 (1.75) | ✓ |
| Bar | 0.606 | | |
| Pie | **0.621** | 2.78 (1.81) | ✓ |
| Bar | 0.580 | | |

stable across visualisation types, the recallability of pie plots is the best in F-questions and RV-questions, bar plots in T-questions and FE-questions, and line plots in U-questions.

We also conducted a user study (N = 8) with experienced scientific researchers who are accustomed to creating their own data visualisations. The four bar plot-pie plot pairs created in the previous crowd-sourced study were presented to the study participants, and they were asked to provide their preferred visualisation for each pair using a 5-point Likert scale (with 1 = the first visualisation is better and 5 = the second visualisation is better (see Figure 8)). We predicted the overall recallability scores for each visualisation using RecallNet and compared it to the preference of our study participants (see Table 4, and supplementary materials for full table). For three of the four visualisation pairs, the ranking of recallability scores agreed with crowd-sourced data. In the only contradictory pair (see the right of Figure 8), study participants opted for the pie visualisation type with a mean score of 3.11 (SD = 1.73).

## 6 DISCUSSION

This work made a substantial leap towards quantifying fine-grained recallability scores on information visualisations.

**VisQA Dataset.** VisQA is the first dataset to introduce fine-grained recallability on an information visualisation dataset as well as high-quality question-answer annotations. The recallability scores are metrics that reveal human performance with a specific type of question. With rich annotations of the elements necessary for the answers, the recallability score of a certain question could be converted into 2D spatial representations (e.g. recallability heatmaps). Since a better visual encoder benefits VQA models [11], the recallability maps could be introduced as an additional

input to VQA models. Addtionally, VisQA is a novel visualisation VQA dataset that uses real-world, visually rich visualisations coming in part from the MASSVIS dataset. Crowd-sourcing is the standard approach for collecting questions on VQA datasets, and the questions for current visualisation VQA datasets [38, 39] were collected by regular crowd workers. In contrast, all the questions in our VisQA came from visualisation experts, which promises a higher quality of questions than previous VQA datasets. Moreover, most visualisations in current VQA datasets [38] are generated pragmatically. However, when it comes to real-world visualisations, the vector representations are usually missing, and researchers have to retrieve the structural information, often by manual annotation [7], which is time-consuming and constrains the dataset size. The introduction of recallability to the VQA setting and the high quality of visualisations and questions enable VisQA to trigger fundamental studies on visualisation QA or chart QA.

**Recallability vs. Recognisability (Memorability).** The bottom-right quadrant in Figure 6 (Right) is completely empty, which means that there are no such visualisations with high recallability (top 15 %) and low memorability (bottom 15 %) in VisQA. This suggests that *memorability is the basis for recallability*, and that *visualisations have to be sufficiently memorable before they become recallable*. The visualisations in the top-right quadrant share some characteristics, like a big

Fig. 8: Survey interface for evaluating visualisation bar-pie pairs. Visualisations in each pair have the same data sources, colouring scheme, and font attributes.

and highlighted title and some explanatory text. Meanwhile, the visualisations in the top-left quadrant of Figure 6 (Right) have high recognisability and low recallability. Compared to the top-right quadrant, visualisations in the top-left quadrant are less recallable. All visualisations in the top-left quadrant are simple monotone plots with few embellishment (e.g. isotype plots). The visualisations in the bottom-left quadrant are easily forgettable and hard to recall. These visualisations are usually overly complex and don't have meaningful titles or additional explanatory text to convey key messages. Compared to the bottom-left quadrant, all the visualisations in top-left and top-right quadrant are much simpler (low data-to-ink ratios), and always with titles, which aligns well with the findings in previous studies [6, 7]. Therefore, our study on VisQA validated previous results and provided interesting insights into how recallability and recognisability (memorability) are different and connected.

**Visualisation Type Recommendation.** In Table 3, T-question and RV-question recallability scores are stable across types (within 4 % of variation). It suggests that the recallability of T-questions and RV-questions are almost irrelevant with visualisation types. As long as the answer elements are presented in visualisations with similar visual

attributes, such as colouring scheme, font family, font size, and element location, the difficulty of recalling the answers to these questions should be on the same level. As for each visualisation type, the recallability of line plots is the best for U-questions, since they are usually the best choice for interpreting time series data [49]. As for FE- and RV-questions, the recallability of line plots is in the last place. Since readers have to go through multiple elements that are far away from each other to find the answers for FE- and RV-questions [9], line plots are not ideal for these kinds of questions. Table 4 demonstrated that overall recallability was in agreement with the preference of our study participants.

We also observed that the FE-, F-, RV-, and U-question recallability of the triplet study was significantly higher than VisQA, and the T-question recallability was not statistically significantly higher. The triplet study contained simple scientific plots, which are further away from real-world, visually rich visualisations, while the length of each HIT in the triplet study was only 1/5 of the study of VisQA. The reason for the significant change of recallability in the above four question types could be explained by subjectively easier questions in the study. Alternatively, crowd workers performed much better in the first several questions

compared to the last ones. To validate this hypothesis, we calculated the question type recallability on VisQA that appeared in the first 1/5 of each HIT. The mean recallability of FE-questions increased from 0.43 to 0.46, F-questions from 0.46 to 0.52, RV-questions from 0.40 to 0.42, and U-questions from 0.48 to 0.42. None of these recallability scores was higher than any recallability score in the triplet study, so the length of the study was not the critical reason for the significant recallability changes. In conclusion, RecallNet generally agreed with the preference of scientific researchers in the use case of visualisation type recommendations.

**Limitations and Future Work.** There is always a trade-off between quality and quantity, which was also the case when designing and collecting our VisQA dataset. Due to the increasing workload in designing high-quality questions for the VQA settings that were specifically targeted for each visualisation, the scale of VisQA became relative small. This influences some computational models, e.g. it caused the over-fitting problem of UMSI [45]. To allow more complex models for recallability prediction, it is essential to extend our VisQA. In the future, we plan to enrich it with more complex data visualisations such as box plots, radar and combination plots. On the other hand, gaze behaviour analysis in a VQA setting on information visualisations has not yet been studied. However, it is a fundamental step to understand the human visual attention system while viewing visualisations. While physical laboratory studies require special-purpose eye tracking equipment, online crowd-sourcing studies or gaze estimation from substitution devices (e.g., mouse, web camera) can be used as a proxy to human attention. In the future, we will investigate such methods to collect human attention data and extend VisQA with such annotations.

## 7 CONCLUSION

This work presented a novel adaptation of a VQA-based study to collect VisQA, a novel visualisation VQA dataset with 200 "in-the-wild" visualisations annotated with crowd-sourced human recallablity scores in five question types, along with a deep convolutional network to predict fine-grained recallability of visualisations. This work made a substantial leap towards quantifying fine-grained recallability scores on information visualisations and envisions several potential applications. We prototypically demonstrated one application developed out of this work. Through a user study, we demonstrated that the prediction from our RecallNet not only maximised user recallability but also agreed with the preferences from scientific researchers in three out of four visualisation triplets, i.e. different visualisations for the same data source. For visualisation type recommendation systems, leveraging recallability would be a strong criterion in providing feedback to users.

## REFERENCES

[1] W. A. Bainbridge, D. D. Dilks, and A. Oliva, "Memorability: A stimulus-driven perceptual neural signature distinctive from memory," *NeuroImage*, vol. 149, pp. 141–152, 2017.

[2] Z. Bylinskii, L. Goetschalckx, A. Newman, and A. Oliva, "Memorability: An image-computable measure of information utility," *arXiv preprint arXiv:2104.00805*, 2021.

[3] L. Standing, "Learning 10000 pictures," *The Quarterly journal of experimental psychology*, vol. 25, no. 2, pp. 207–222, 1973.

[4] P. Isola, J. Xiao, A. Torralba, and A. Oliva, "What makes an image memorable?" in *CVPR 2011*. IEEE, 2011, pp. 145–152.

[5] A. Khosla, A. S. Raju, A. Torralba, and A. Oliva, "Understanding and predicting image memorability at a large scale," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 2390–2398.

[6] M. A. Borkin, A. A. Vo, Z. Bylinskii, P. Isola, S. Sunkavalli, A. Oliva, and H. Pfister, "What makes a visualization memorable?" *IEEE transactions on visualization and computer graphics*, vol. 19, no. 12, pp. 2306–2315, 2013.

[7] M. A. Borkin, Z. Bylinskii, N. W. Kim, C. M. Bainbridge, C. S. Yeh, D. Borkin, H. Pfister, and A. Oliva, "Beyond memorability: Visualization recognition and recall," *IEEE transactions on visualization and computer graphics*, vol. 22, no. 1, pp. 519–528, 2015.

[8] N. C. Rust and V. Mehrpour, "Understanding image memorability," *Trends in cognitive sciences*, vol. 24, no. 7, pp. 557–568, 2020.

[9] P. Polatsek, M. Waldner, I. Viola, P. Kapec, and W. Benesova, "Exploring visual attention and saliency modeling for task-based visual analysis," *Computers & Graphics*, vol. 72, pp. 26–38, 2018.

[10] S. Antol, A. Agrawal, J. Lu, M. Mitchell, D. Batra, C. L. Zitnick, and D. Parikh, "Vqa: Visual question answering," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 2425–2433.

[11] Y. Goyal, T. Khot, D. Summers-Stay, D. Batra, and D. Parikh, "Making the v in vqa matter: Elevating the role of image understanding in visual question answering," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 6904–6913.

[12] S. E. Kahou, V. Michalski, A. Atkinson, Á. Kádár, A. Trischler, and Y. Bengio, "Figureqa: An annotated figure dataset for visual reasoning," *arXiv preprint arXiv:1710.07300*, 2017.

[13] K. Kafle, R. Shrestha, S. Cohen, B. Price, and 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.

[14] R. Amar, J. Eagan, and J. Stasko, "Low-level components of analytic activity in information visualization," in *IEEE Symposium on Information Visualization, 2005. INFOVIS 2005.* IEEE, 2005, pp. 111–117.

[15] T. F. Brady, T. Konkle, G. A. Alvarez, and A. Oliva, "Visual long-term memory has a massive storage capacity for object details," *Proceedings of the National Academy of Sciences*, vol. 105, no. 38, pp. 14 325–14 329, 2008.

[16] P. Isola, D. Parikh, A. Torralba, and A. Oliva, "Understanding the intrinsic memorability of images," MASSACHUSETTS INST OF TECH CAMBRIDGE, Tech. Rep., 2011.

[17] P. Isola, J. Xiao, D. Parikh, A. Torralba, and A. Oliva, "What makes a photograph memorable?" *IEEE transactions on pattern analysis and machine intelligence*, vol. 36, no. 7, pp. 1469–1482, 2013.

[18] M. Mancas and O. Le Meur, "Memorability of natural scenes: The role of attention," in *2013 IEEE International Conference on Image Processing.* IEEE, 2013, pp. 196–200.

[19] S. Perera, A. Tal, and L. Zelnik-Manor, "Is image memorability prediction solved?" in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, 2019, pp. 0–0.

[20] A. Jaegle, V. Mehrpour, Y. Mohsenzadeh, T. Meyer, A. Oliva, and N. Rust, "Population response magnitude variation in inferotemporal cortex predicts image memorability," *Elife*, vol. 8, p. e47596, 2019.

[21] J. Fajtl, V. Argyriou, D. Monekosso, and P. Remagnino, "Amnet: Memorability estimation with attention," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 6363–6372.

[22] F. Haist, A. P. Shimamura, and L. R. Squire, "On the relationship between recall and recognition memory." *Journal of Experimental Psychology: Learning, Memory, and Cognition*, vol. 18, no. 4, p. 691, 1992.

[23] A. P. Yonelinas, "The nature of recollection and familiarity: A review of 30 years of research," *Journal of memory and language*, vol. 46, no. 3, pp. 441–517, 2002.

[24] J. Hullman, M. Kay, Y.-S. Kim, and S. Shrestha, "Imagining replications: Graphical prediction & discrete visualizations improve recall & estimation of effect uncertainty," *IEEE transactions on visualization and computer graphics*, vol. 24, no. 1, pp. 446–456, 2017.

[25] W. A. Bainbridge, E. H. Hall, and C. I. Baker, "Drawings of real-world scenes during free recall reveal detailed object and spatial information in memory," *Nature communications*, vol. 10, no. 1, pp. 1–13, 2019.

[26] W. S. Cleveland and R. McGill, "Graphical perception: Theory, experimentation, and application to the development of graphical methods," *Journal of the American statistical association*, vol. 79, no. 387, pp. 531–554, 1984.

[27] S. M. Kosslyn, "Understanding charts and graphs," *Applied cognitive psychology*, vol. 3, no. 3, pp. 185–225, 1989.

[28] S. Pinker, "A theory of graph comprehension," *Artificial intelligence and the future of testing*, pp. 73–126, 1990.

[29] S. Bateman, R. L. Mandryk, C. Gutwin, A. Genest, D. McDine, and C. Brooks, "Useful junk? the effects of visual embellishment on comprehension and memorability of charts," in *Proceedings of the SIGCHI conference on human factors in computing systems*, 2010, pp. 2573–2582.

[30] R. Borgo, A. Abdul-Rahman, F. Mohamed, P. W. Grant, I. Reppa, L. Floridi, and M. Chen, "An empirical study on using visual embellishments in visualization," *IEEE Transactions on Visualization and Computer Graphics*, vol. 18, no. 12, pp. 2759–2768, 2012.

[31] A. V. Moere, M. Tomitsch, C. Wimmer, B. Christoph, and T. Grechenig, "Evaluating the effect of style in information visualization," *IEEE transactions on visualization and computer graphics*, vol. 18, no. 12, pp. 2739–2748, 2012.

[32] X. Shu, A. Wu, J. Tang, B. Bach, Y. Wu, and H. Qu, "What makes a data-gif understandable?" *IEEE Transactions on Visualization and Computer Graphics*, vol. 27, no. 2, pp. 1492–1502, 2020.

[33] S.-H. Kim, Z. Dong, H. Xian, B. Upatising, and J. S. Yi, "Does an eye tracker tell the truth about visualizations?: Findings while investigating visualizations for decision making," *IEEE Transactions on Visualization and Computer Graphics*, vol. 18, no. 12, pp. 2421–2430, 2012.

[34] Y.-S. Kim, K. Reinecke, and J. Hullman, "Explaining the gap: Visualizing one's predictions improves recall and comprehension of data," in *Proceedings of the 2017 CHI Conference on Human Factors in Computing Systems*, 2017, pp. 1375–1386.

[35] H.-K. Kong, Z. Liu, and K. Karahalios, "Trust and recall of information across varying degrees of title-visualization misalignment," in *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*, 2019, pp. 1–13.

[36] X. Chen, M. Jiang, and Q. Zhao, "Predicting human scanpaths in visual question answering," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 10 876–10 885.

[37] K. Kafle, B. Price, S. Cohen, and 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.

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

[39] R. Chaudhry, S. Shekhar, U. Gupta, P. Maneriker, P. Bansal, and 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.

[40] A. Newman, B. Mcnamara, C. Fosco, Y. Zhang, P. Sukhum, M. Tancik, N. Kim, and Z. Bylinskii, "Turkeyes: A web-based toolbox for crowdsourcing attention data," in *CHI '20: CHI Conference on Human Factors in Computing Systems*, 2020.

[41] F. Chollet, "Xception: Deep learning with depthwise separable convolutions," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 1251–1258.

[42] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille, "Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution,

and fully connected crfs," *IEEE transactions on pattern analysis and machine intelligence*, vol. 40, no. 4, pp. 834–848, 2017.

[43] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.

[44] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.

[45] C. Fosco, V. Casser, A. K. Bedi, P. O'Donovan, A. Hertzmann, and Z. Bylinskii, "Predicting visual importance across graphic design types," in *Proceedings of the 33rd Annual ACM Symposium on User Interface Software and Technology*, 2020, pp. 249–260.

[46] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *2009 IEEE conference on computer vision and pattern recognition*. Ieee, 2009, pp. 248–255.

[47] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.

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

[49] Y. Wang, F. Han, L. Zhu, O. Deussen, and B. Chen, "Line graph or scatter plot? automatic selection of methods for visualizing trends in time series," *IEEE transactions on visualization and computer graphics*, vol. 24, no. 2, pp. 1141–1154, 2017.

**Mihai Bâce** is a post-doctoral researcher in the Perceptual User Interfaces group at the University of Stuttgart, Germany. He did his PhD at ETH Zurich, Switzerland, at the Institute for Intelligent Interactive Systems. He received his MSc. in Computer Science from École Polytechnique Fédérale de Lausanne, Switzerland, and his BSc. in Computer Science from the Technical University of Cluj-Napoca, Romania. His research interests include computational Human-Computer Interaction with a focus on sensing and modelling user attention.

**Andreas Bulling** is Full Professor of Computer Science at the University of Stuttgart, Germany, where he directs the research group "Human-Computer Interaction and Cognitive Systems". He received his MSc. in Computer Science from the Karlsruhe Institute of Technology, Germany, in 2006 and his PhD in Information Technology and Electrical Engineering from ETH Zurich, Switzerland, in 2010. Before, Andreas Bulling was a Feodor Lynen and Marie Curie Research Fellow at the University of Cambridge, UK, and a Senior Researcher at the Max Planck Institute for Informatics, Germany. His research interests include computer vision, machine learning, and human-computer interaction.

**Yao Wang** is a Ph.D. student at the University of Stuttgart, Germany. He received his BSc. in Intelligence Science and Technology and MSc. in Computer Software and Theory both from Peking University, China, in 2017 and 2020, respectively. His academic research interest focuses on visual attention modelling on information visualizations.

**Chuhan Jiao** is a Master's student in Computer Science at Aalto University. He received his BEng. in Computer Science and Technology from Donghua University, China. His research interest lies at the intersection of Human-Computer Interaction and Computer Vision.

# Supplementary Material:
# VisQA: Quantifying information Visualisation Recallability via Question Answering

Yao Wang, Chuhan Jiao, Mihai Bâce, and Andreas Bulling

**Figure II-2: 2011 Critical Technology Transactions by Region of Foreign Acquirer**

Question: What is the theme of this visualization?
**A: Critical Technology Transactions by Region of Foreign Acquirer**
B: Critical Product Transactions by Region of Foreign Acquirer
C: Critical Project Transactions by Region of Foreign Acquirer
D: I can not remember
Type: T-question

Question: What year's data is displayed in this visualization?
**A: 2011**
B: 2012
C: 2013
D: I can not remember
Type: F-question

Question: Which region has the lowest value?
A: Middle East & North Africa
**B: Other**
C: East Asia
D: I can not remember
Type: FE-question

Question: Which region has the highest value?
A: East Aisa
B: Canada, Australia & New Zealand
**C: Europe (excluding Russia)**
D: I can not remember
Type: FE-question

Question: How many transactions does Europe(excluding Russia) have?
A: Around 20
B: Around 50
**C: Around 70**
D: I can not remember
Type: RV-question

Fig. 1: Example visualisation of *bar plot* from VisQA dataset with five multiple-choice questions.

Question: What is the theme of this visualization?
**A: Reno housing unit growth outpaced population and house hold growth during the past decade**
B: Reno population growth outpaced housing unit and house hold growth during the past decade
C: Reno house hold growth outpaced housing unit and population growth during the past decade
D: I can not remember
Type: T-question

Question: What period of data does this visualization show?
A: 2010-2020
**B: 2000-2010**
C: 2000-1990
D: I can not remember
Type: F-question

Question: Which one has the highest annual growth rate?
A: Reno-Sparks Population
B: Reno-Sparks Households
**C: Reno-Sparks Housing Units**
D: I can not remember
Type: FE-question

Question: Which one is higher, annual growth rate of Reno-Sparks Households or annual growth rate of Reno-Sparks Housing Units?
**A: Reno-Sparks Housing Units**
B: Reno-Sparks Households
C: The same
D: I can not remember
Type: U-question

Question: What is the population in 2000
**A: Around 350,000**
B: Around 400,000
C: Around 450,000
D: I can not remember
Type: RV-questio

Fig. 2: Example visualisation of *table* from VisQA dataset with five multiple-choice questions.



Question: What is the theme of this visualization?
A: Tourism and GDP
B: Foreign exchange reserves and GDP
**C: Military spending and GDP**
D: I can not remember
Type: T-question

Question: What period of data does this visualization show?
A: 2001-2010
B: 2002-2010
**C: 2002-2011**
D: I can not remember
Type: F-question

Question: Which country has the biggest GDP growth in this period?
A: China
**B: Angola**
C: Ethiopia
D: I can not remember
Type: FE-question

Question: Which country has the biggest Military spending growth in this period?
A: Ecuador
**B: Kazakhstan**
C: Armenia
D: I can not remember
Type: FE-question

Question: Which country has a negative growth in military spending?
**A: Italy**
B: Canada
C: Singapore
D: I can not remember
Type: F-question

Fig. 3: Example visualisation of *scatter plot* from VisQA dataset with five multiple-choice questions.

**World imports**
% of total

Developed economies          F'CAST

70
60
50
40
Emerging economies
30
20
10
0

1990 91 92 93 94 95 96 97 98 99 2000 01 02 03 04 05 06 07 08 09 10 11 12
Sources: WTO; The Economist

Question: What is the theme of this visualization?
A: World economies
**B: World imports**
C: World exports
D: I can not remember
Type: T-question

Question: What period of data does this visualization show?
A: 1990-2010
B: 1990-2011
**C: 1990-2012**
D: I can not remember
Type: F-question

Question: What range of percentage does this visualization show?
A: 0-60
**B: 0-70**
C: 0-80
D: I can not remember
Type: RV-question

Question: What decreases as time goes by?
**A: Developed economies**
B: Emerging economies
C: General economies
D: I can not remember
Type: U-question

Question: What increases as time goes by?
A: Developed economies
**B: Emerging economies**
C: General economies
D: I can not remember
Type: U-question

Fig. 4: Example visualisation of *line plot* from VisQA dataset with five multiple-choice questions.



**Total Professional, Scientific, and Technical Services**
**2009-2011**

Other Professional, Scientific, and Technical Services
4%
Scientific Research and Development Services
10%
Management, Scientific, and Technical Consulting Services
15%
Computer Systems Design and Related Services
44%
Architectural, Engineering, and Related Services
27%

**Covered Transactions from the Professional, Scientific, and Technical Services Subsector**

Question: What is the theme of this visualization?
A: Total Professional services
B: Total Professional services and Scientific services
**C: Total Professional services, Scientific and technical services**
D: I can not remember
Type: T-question

Question: What year's data is displayed in this visualization?
A: 2009-2010
**B: 2009-2011**
C: 2007-2009
D: I can not remember
Type: F-question

Question: Which Field has the highest share of contribution?
**A: Computer systems and related services**
B: Scientific research and development
C: Architectural Engineering
D: I can not remember
Type: FE-question

Question: Which Field has the lowest share of contribution?
A: Computer systems and related services
B: Scientific research and development
**C: Other professional, scientific and technical services**
D: I can not remember
Type: FE-question

Question: Which keyword stands for the visualization the best?
A: Professional
**B: Transactions**
C: Technologies
D: I can not remember
Type: F-question

Fig. 5: Example visualisation of *pie chart* from VisQA dataset with five multiple-choice questions.

Question: What is the theme of this visualization?
**A: Beer consumption around the world**
B: Beer consumption around Europe only
C: Beer consumption around US only
D: I can not remember
Type: T-question

Question: What are the years seen in the visualization?
A: 2008 and 2010
B: 2000 and 2001
**C: 2008 and 2009**
D: I can not remember
Type: F-question

Question: Which country has recorded the highest consumption in terms of literes per person?
**A: Czech Republic**
B: US
C: Germany
D: I can not remember
Type: F-question

Question: Which country has recorded the lowest consumption in terms of literes per person?
**A: Vietnam**
B: France
C: China
D: I can not remember
Type: F-question

Question: How much Hectoliters/m of beer is produced by Asia and Europe in the year 2009?
A: Around 500
**B: Above 550**
C: Below 400
D: I can not remember
Type: RV-question

Fig. 6: Example visualisation of *other* types from VisQA dataset with five multiple-choice questions.



Question: What is the theme of this visualization?
**A: Number of vaccinated white people according to age group**
B: Number of vaccinated black people according to age group
C: Number of vaccinated asian people according to age group
D: I can not remember
Type: T-question

Question: Which age group is vaccinated to the maximum?
A: 65-69
**B: 70-79**
C: 55-59
D: I can not remember
Type: FE-question

Question: What is the mimumum value of vaccinated people?
**A: Less than one million**
B: More than one million
C: Equal to one million
D: I can not remember
Type: RV-question

Question: What is the color used for of denoting the maximum?
**A: dark blue**
B: light blue
C: light pink
D: I can not remember
Type: U-question

Question: Which country's data is used in the visualization?
**A: UK**
B: US
C: Europe
D: I can not remember
Type: F-question

Fig. 7: Triplet visualisations (pie, line, bar) created by the same data source (Group1)

Question: What is the theme of this visualization?
**A: Average Cohort Outcomes in New York City**
B: Average Cohort Outcomes in Washington
C: Average Cohort Outcomes in Los Angeles
D: I can not remember
Type: T-question

Question: Which school had the maximum average of cohorts?
A: PACE HIGH SCHOOL
**B: LOWER EAST SIDE PREPARATORY HIGH SCHOOL**
C: CASCADES HIGH SCHOOL
D: I can not remember
Type: FE-question

Question: What is the source of the visualization?
**A: New York City Department of Education data**
B: Washington Department of Education data
C: Los Angeles Department Education Data
D: I can not remember
Type: F-question

Question: What is the maximum average of cohorts reported in the visualization?
**A: More than 200**
B: Less than 200
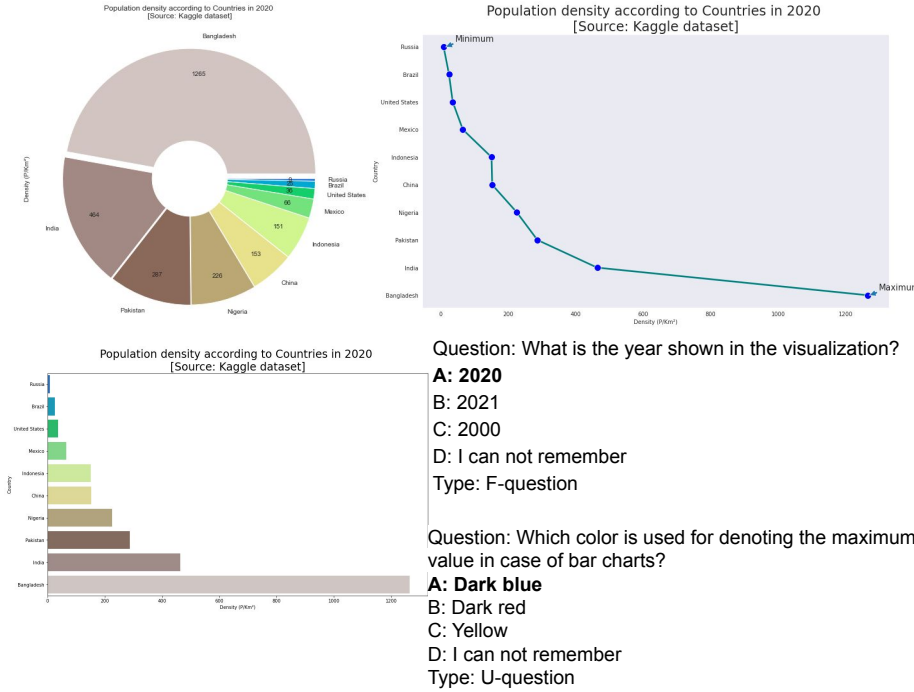C: Equal to 200
D: I can not remember
Type: RV-question

Question: In bar charts, Which color is used for denoting the mimimum in all cases?
**A: Shades of blue**
B: Shades of red
C: Shades of yellow
D: I can not remember
Type: U-question

Fig. 8: Triplet visualisations (pie, line, bar) created by the same data source (Group2)

TABLE 1: Full questionnaire results. Likert-scores are presented in mean and standard deviation. 1 = strongly agreed, and 5 = strongly disagreed.

| Question | Percentage or Likert-score |
|---|---|
| Please select your gender | Male(56%) Female(33%) N/A(11%) |
| How often do you need to create your own visualisation? | About once per week(33%) About once per month(33%) Two or more times per month(33%) |
| Please select the tools you have used before. [Office (Word, Excel, PPT...)] | Yes(78%) No(22%) |
| Please select the tools you have used before. [R] | Yes(56%) No(44%) |
| Please select the tools you have used before. [python] | Yes(78%) No(22%) |
| Please select the tools you have used before. [matlab] | Yes(89%) No(11%) |
| Sometimes I don't know which visualisation type is the best for my data. | $3.00 \pm 0.87$ |
| I have a strong preference for a plot type, and I will try to use that type as much as possible. | $2.78 \pm 1.20$ |
| It would be useful to know how effective my visualisation is to convey information to readers. | $4.89 \pm 0.33$ |
| If there is a tool to automatically help me to decide the visualisation type (just like the previous decisions between bar and pie plots), I would like to try it out. | $4.33 \pm 0.87$ |

Average age of survivors in Titanic according to class

Question: What is the theme of this visualization?
A: Average age of survivors in Titanic
**B: Average age of youth survivors in Titanic**
C: Average age of children survivors in Titanic
D: I can not remember
Type: T-question

Question: Which category had the maximum average of survivors?
A: Second
**B: First**
C: Third
D: I can not remember
Type: FE-question

Question: What is the source of the visualization?
**A: The Encyclopedia Titanica**
B: The Kaggle Titanica
C: Titanica
D: I can not remember
Type: F-question

Question: What is the minimum average age reported in the visualization?
**A: More than 20**
B: Less than 20
C: Equal to 20
D: I can not remember
Type: RV-question

Question: Which color pair is used for of denoting the survivors?
**A: Blue-Green**
B: Red-Pink
C: Yello-Pink
D: I can not remember
Type: U-question

Fig. 9: Triplet visualisations (pie, line, bar) created by the same data source (Group3)

Fig. 10: Triplet visualisations (pie, line, bar) created by the same data source (Group4)

| Visualisations | Memorability (Recognisability) | Recallability |
|---|---|---|
|  | 3.986 | 0.358 |
|  | 3.669 | 0.183 |
|  | 3.463 | 0.075 |
|  | 3.382 | 0.292 |

TABLE 2: Full memorability (recognisability) and recallability scores of all visualisations in top-left quadrant in Figure 6 (Right) from the main manuscript.
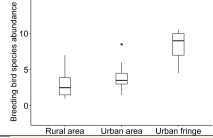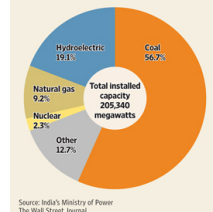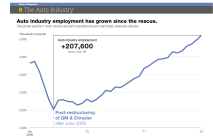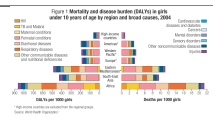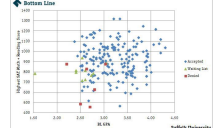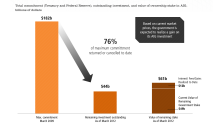
| Visualisations | Memorability (Recognisability) | Recallability |
|---|---|---|
|  | 4.268 | 0.633 |
|  | 4.203 | 0.679 |
|  | 3.986 | 0.658 |
|  | 3.986 | 0.650 |

TABLE 3: Full memorability (recognisability) and recallability scores of all visualisations in top-right quadrant in Figure 6 (Right) from the main manuscript.
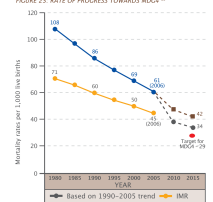
| Visualisations | Memorability (Recognisability) | Recallability |
|---|---|---|
|  | 0.459 | 0.318 |
|  | 0.855 | 0.317 |
|  | 1.139 | 0.337 |
|  | 1.055 | 0.323 |

TABLE 4: Full memorability (recognisability) and recallability scores of all visualisations in bottom-left quadrant in Figure 6 (Right) from the main manuscript.