# Computational Perception for Multi-Modal Document Understanding

by

Zoya Bylinskii

B.S., Computer Science, University of Toronto, 2012
S.M., Electrical Engineering and Computer Science, M.I.T., 2015

Submitted to the Department of Electrical Engineering and Computer Science
in partial fulfillment of the requirements for the degree of

Doctor of Philosophy
in Electrical Engineering and Computer Science
at the Massachusetts Institute of Technology

September 2018

## Computational Perception for Multimodal Document Understanding
by Zoya Bylinskii

Submitted to the Department of Electrical Engineering and Computer Science
on August 31, 2018, in partial fulfillment of the requirements for the degree of
Doctor of Philosophy

**Abstract**

Multimodal documents occur in a variety of forms, as graphs in technical reports, diagrams in textbooks, and graphic designs in bulletins. Humans can efficiently process the visual and textual information contained within to make decisions on topics including business, healthcare, and science. Building the computational tools to understand multimodal documents can have important applications for web search, information retrieval, captioning and summarization, and automated design. This thesis makes contributions on two fronts: (i) to the development of data collection methods for measuring how humans perceive multimodal documents (i.e., where they look, what they find important), and (ii) to the development of computer vision tools for automatically parsing and making predictions about multimodal documents (i.e., the subject matter they are about). Specifically, the crowdsourced attention data captured from our novel user interfaces is used to train neural network models to predict where people look in graphic designs and information visualizations, with demonstrated applications to thumbnailing, design retargeting, and interactive feedback within graphic design tools. Separately, our models for detecting visual elements and parsing text elements in infographics (information graphics) are used for topic prediction and to present a system for automatic summarization. This thesis makes contributions at the interface of human and computer vision, with applications to human-computer interfaces and design.

Thesis Supervisor: Fredo Durand
Title: Professor of Electrical Engineering and Computer Science

Thesis Supervisor: Aude Oliva
Title: Principal Research Scientist, MIT CSAIL

# Acknowledgments

**Thesis committee:** The make-up of my committee is a testament to the interdisciplinary nature of this thesis, as the intersection of 5 fields: Fredo Durand (computer graphics), Aude Oliva (cognitive science, perception), Hanspeter Pfister (visualization), Rob Miller (human-computer interaction), and Bill Freeman (computer vision). I could not have chosen a better set of researchers to represent these 5 fields and serve as my role models within them. I thank these academics for the influences they have had on my academic work and papers, through their own examples, their general advice to graduate students, and through direct feedback. *Novelty, clarity, impact,* and *thoroughness* are probably the adjectives most representative of the approach to research they follow. I have borrowed at least a little bit of wisdom from each of them, and their cumulative wisdom has set me on the path to continuing my research career. The interdisciplinary path is not an easy one, often lacking a stable community and clear rules to play by, and requiring bouts of swimming against the current; however, it is a very rewarding path that naturally opens up many opportunities for novel and impactful work.

**Advisors:** The burden on academic advisors to support students, financially, academically, and emotionally, for the duration of a 4-6 year PhD program, is an immense one. I am very thankful to Fredo Durand and Aude Oliva for their incredible support throughout my PhD journey, for always finding the time to meet, and for buckets of advice and encouragement. Fredo has taught me the necessity of precision - in every sentence, word choice, and discussion; thoroughness in tackling research problems and understanding results; and having explanations for observations at the micro-level. Aude has taught me to be *"crisp", "punchy",* and *"memorable"* - both in the research ideas I choose to pursue and how I communicate the final results. She has helped me cultivate both research and leadership skills, and has supported every initiative I've taken on. Most importantly, she has taught me to *"think decades ahead".*

**Mentors:** I have been incredibly lucky in the amount of mentors that made their time and advice available to me whenever I needed it. I thank Sven Dickinson and John Tsotsos for staying in touch with me throughout my PhD and offering their unwavering support. I thank Aaron Hertzmann and Bryan Russell for their wonderful advising during my Adobe research internship - a treasure trove of research input I wish I would have discovered even earlier in my PhD. I thank Stefanie Mueller for selflessly donating

her time to offer career advice and resources. I thank also the many senior graduate students like Adrian Dalca who were always willing to share their experiences, help problem solve, and provide encouragement.

**Colleagues:** I began my PhD working with senior students from different fields, most notably Melissa Le-Hoa Võ (human perception), Phillip Isola (computer vision), and Michelle Borkin (visualization). All three, who were amazing mentors as graduate students and postdocs, are now faculty. My first papers were due to these collaborations, and thanks to the wisdom, ideas, and technical expertise of these individuals, who have taught me a great deal about the methods and expectations in their respective fields. I thank also my frequent collaborators Nam Wook Kim (visualization, human-computer interaction) and Adrià Recasens (computer vision) for their insights, efficiency, and enthusiasm to take on novel challenges. I also immensely enjoyed collaborations and interactions with Tilke Judd (saliency) and Peter O'Donovan (design tools).

Progressing through my PhD, I learned how rewarding mentoring and advising younger students can be, and have over time grown the group of students I work with. I particularly want to thank Spandan Madan, Sami Alsheikh, Matthew Tancik, Anelise Newman, Kimberli Zhong, Nathan Landman, and Camilo Fosco for being part of my *"vizteam"* over the last two years of my PhD and for their incredible research drives that have led to multiple submitted manuscripts, abstracts, and poster presentations. I am confident in the research potential of these students and look forward to seeing their next steps.

**Family and friends:** The unconditional support and belief in one's success provided by family and friends is often enough to help one sail through a PhD. The encouragement of my parents to learn the tools of computer science, and the advice of my friends and mentors through undergrad to pursue research opportunities left no question in my mind that I was going to complete a PhD and continue along the research path. My dad's words throughout my childhood that *"a minute not spent learning is a minute lost"* and *"an easy path is not a path worth taking"* have burned into my conscience and influenced the decisions that got me to this point. The people who have cultivated my interests and have had the largest impacts on my education and outlook on the world are also the people who I continue to look up to most: my parents, Marina & Vasili Gavrilov and my husband Alexei Bylinskii.

**Crowdworkers:** My thesis would not have happened without the efforts of hundreds of crowdworkers all over the world, willing to click on, stare at, and label thousands of images. I am incredibly thankful to all the anonymous participants on Amazon's Mechanical Turk for their input into this thesis. When online experiments were not suitable for one reason or another, CSAIL and the broader MIT community came to the rescue. Dozens of MIT students were willing to donate their time in exchange for money, and sometimes even just burritos, and sometimes with less than an hour's notice (Fig. 1). Thank you to my *"vizteam"* students for helping set up such efficient data collection pipelines!

**Figure 1.** A lot of the work in this thesis has been built on top of the crowdsourcing efforts of participants both online and at MIT. I would like to thank the hundreds of participants that provided the annotation data used for training and evaluating our computational models (in exchange for money or just burritos). Thank you to my *"vizteam"* students for helping set up such efficient data collection pipelines!

# Contents

## II State of Computational Models of Attention 53

## 6 Where should saliency models look next? 55

## III Crowdsourcing Human Attention 71

## 7 BubbleView: an interface for crowdsourcing image importance 73

## 8 ZoomMap: using zoom to capture user areas of interest on images 111

# List of Figures

# Chapter 1

# Introduction

*If an image is worth a thousand words, then a multimodal document is worth a thousand concepts.*

**H**UMANS can efficiently learn from multimodal (e.g., visual and textual) content present in a diverse set of sources such as textbooks, slideshow presentations, and posters (Fig. 1.1). We regularly extract information from illustrations in textbooks, parse graphs and charts to make decisions, and allow informational posters to influence our opinions on different topics, as the visuals burn into our memory. High-level cognitive abilities are required to integrate information from the textual and visual modalities, to reason about the structure of the documents (e.g., how the accompanying figures support the text), and to summarize the most important points for decision-making.

**How do humans do this?** One way to answer this question is by studying the

**Figure 1.1.** Examples of multimodal documents: (a) information visualizations, (b) graphic designs (posters), (c) articles and documents with figures, (d) presentation slides, (e) infographics (information graphics). See Section 1.1 for a taxonomy of multimodal documents.

human perception of multimodal documents: what catches people's attention, what they spend time studying, and what they ultimately find most important. The first parts of this thesis (II-IV) are about human attention, how to crowdsource, analyze, and model it, with a specific focus on multimodal documents, including information visualizations (Sec. 7.3.1, 9.3.1), webpages (Sec. 7.3.3), graphic designs (Sec. 7.4.1, 9.3.2), and academic posters (Sec. 8.4). Novel data collection methods for crowdsourcing human attention at-scale are introduced. Through controlled lab experiments, we show that eye movements collected using a standard eye tracker can be approximated by (i) having users click to view regions of an image that have been initially blurred (Fig. 1.2a; Chapter 7), and (ii) having users zoom in to an image to view regions at larger resolution (Fig. 1.2b; Chapter 8). The collected data can then be used to train computational models of attention and for automatic design applications like retargeting and thumbnailing (Chapters 9-10). An understanding of the human perception of multimodal documents can thus inform computational approaches. Note that Section 1.1 provides an overview of the types of multimodal documents that are the focus of this thesis.

**Figure 1.2.** Two novel data collection methods for crowdsourcing human attention patterns on images: (a) The BubbleView interface presents blurred images to participants, and they can click around to uncover regions of the image ("bubbles") at full resolution. The pattern of clicks, when averaged over multiple participants to produce a click map, approximate an attention map for the image (Chapter 7). (b) The ZoomMap interface allows participants to explore images with the pinch zoom gesture on their mobile devices. The zoom actions can be converted to a heatmap of important regions in an image (Chapter 8).



**Why develop computational approaches?** Automatically parsing multimodal documents including information visualizations, graphic designs, infographics, and presentation slides (see Fig. 1.1 and Sec. 1.1) can make a variety of applications possible. For example, we may want to automatically convert a standard information visualization like a plain bar graph into a more effective and memorable infographic (Fig. 1.3a); within graphic design tools, we may want to provide automatic feedback to a designer

**Figure 1.3.** Examples of applications with multimodal documents: (a) turning standard visualizations into effective infographics, (b) providing automatic feedback within graphic design tools (Chapter 9), (c) summarizing multimodal content (Chapters 8.4, 10, 13), (d) annotating documents with visuals, (e) translating text into visual presentations.



about the parts of their design that are likely to attract observer attention (Fig. 1.3b); given a design, we may also propose visual and textual summaries (Fig. 1.3c); we could learn to automatically annotate text articles with relevant visuals, to help guide the reader's attention and increase reader engagement (Fig. 1.3d); more ambitiously, we could approach the problem of translating the knowledge in text documents and text-books into visual presentations, slides, and posters (Fig. 1.3e). These applications can help democratize graphic design for areas such as data exploration and internet-scale education, for disseminating information in a broadly-accessible format.

**The technical challenges:** To work towards these applications, we need to first develop the tools for parsing multimodal documents. However, multimodal documents pose many technical challenges for computational systems. There is large variability in how the text and visuals are arranged and scaled, the text is rendered in different fonts, and the visual elements range from photo-realistic to abstract icons (Fig. 1.4a). Separate approaches exist for detecting and extracting text from images (i.e., optical character recognition; see Chapter 2), and for parsing the visual content of a scene (e.g., computer vision approaches including object detection, scene classification, and semantic segmentation). However, handling the semantic and stylistic variability in multimodal documents remains a challenge to computational models trained mostly on natural images [95, 98, 218]. Furthermore, integrating information from textual and visual modalities simultaneously and reasoning about document structure is a relatively unexplored research direction at the interface of computer vision, human computer in-

teraction, and natural language processing.  Part V of this thesis contains initial steps in this direction: parsing the text and visual elements in infographics (Fig. 1.4b; Chapter 12) to be able to predict the topics and subtopics being addressed, and to generate multimodal (textual and visual) summaries of the content (Fig. 1.4c; Chapter 13).

**Figure 1.4.** (a) Example infographics, containing text and images with large variations in style, scale, and semantics (Chapter 11).  (b) An infographic that has been automatically parsed.  Annotated are the detected text in green and the detected and classified icons in red (Chapter 12).  (c) A sample multimodal summary automatically produced from the detected text and icons in (b).  The summary captures the main topic of the infographic, and is composed of a representative text tag and visual hashtag (Chapter 13).

The ubiquity of multimodal documents: From business presentations, medical documents, and textbooks to children's books and subway advertisements, we are regularly exposed to multimodal documents. These documents have been curated with a human audience in mind, to effectively convey concepts, deliver messages, and tell stories - in educational, business, and promotional settings. A vast amount of semantic and design knowledge is thus encoded in multimodal documents, and this thesis presents a multi-pronged approach to begin to harness this knowledge: using crowdsourcing techniques to measure how humans perceive multimodal documents (Part III), and simultaneously building the computer vision tools to make predictions about the content in multimodal documents (Parts IV and V).

## 1.1 Taxonomy of multimodal documents

Multimodal documents as referred to in this thesis are intended to cover a broad range of document types, united by the fact that they have been composed out of visual and textual elements. This includes presentation slides, articles with figures, textbooks with diagrams, and information visualizations (Fig. 1.5). Under this category, we also include graphic designs that contain text (e.g., comics, advertisements, and diagrams).

Information visualizations are visual and numerical representations of data, in widely ranging formats - from charts, graphs, and tables, to scientific visualizations (molecular diagrams, medical scans, renderings of astrological phenomena, etc.) and infographics (information graphics). Scientific visualizations and infographics have a lot in common with graphic designs in that they have a strong focus on using visual content (human-recognizable objects, photographs, graphics simulations, and icons/pictographs) for communicative purposes. Fig. 1.6 shows how we think of **infographics**: as lying at the intersection of information visualizations and graphic designs. Where graphic designs are often created to tell a story or deliver a message, and information visualizations are used to depict data, infographics can be seen as *telling a story with data*. Some of the computer vision methods for parsing infographics that will be considered later in this thesis are related also to work on graphic designs without text, like clipart and artistic media.

In this thesis, studies of human perception are carried out on information visualizations, webpages, graphic designs, and academic posters, with a stronger focus on information visualizations and graphic designs (Part III). Computational models of attention are then trained using the human perception data on information visualizations and graphic designs (Part IV). Infographics form a subset of information visualizations. Infographics then become the sole focus of the computer vision tools developed for text and icon detection in the latter part of the thesis (Part V).

## 1.2 Datasets

Natural images: For evaluating computational models of visual attention (Chapter 6), we used the *MIT300* [106] and *CAT2000* [20] datasets from our own MIT

**Figure 1.5.** Multimodal documents are composed of text and images, and include widely ranging media from presentation slides, articles with figures, textbooks with diagrams, and information visualizations. Graphic designs with text, like comics, advertisements, and diagrams are also part of this category. At the intersection of information visualizations and graphic designs lie infographics and scientific visualizations. Infographics, which are the focus of the latter part of this thesis, are information visualizations that tell stories with data and have a lot in common with graphic designs. Related work in computer vision on graphic designs includes papers and datasets on: (1) clipart [232, 233], (2) artistic media [218], (3) comics [98], (4) advertisements [95], and (5) diagrams [109, 193].



Saliency Benchmark [31]. We also used images from the *CAT2000* dataset for evaluating our ZoomMap interface (Chapter 8). We used images from the *OSIE* [223] and *SALICON* [103] datasets for evaluating our BubbleView interface (Chapter 7).

**Visualizations:** We used our own 5K dataset of information visualizations called *MASSVIS*, first introduced in [23, 24] (see Chapter 4). For training a computational model of attention for visualizations, we collected perception data on 1,411 visualizations from this dataset (Chapter 9). These include infographics (from *Visual.ly*) as well as visualizations from news media and government publications. A subset of these visualizations were also used for evaluating our BubbleView interface (Chapter 7).

**Graphic designs:** For training our computational model of attention (Chapter 9), we used the *Graphic Design Importance* (GDI) dataset of 1,078 Flickr graphic designs from [154], provided with importance annotations. A subset of these graphic designs were also used for evaluating our BubbleView interface (Chapter 7). We collected

**Figure 1.6.** Where graphic designs are often created to tell a story or deliver a message, and information visualizations are used to depict data, infographics can be seen as telling a story with data. This is why we place infographics at the intersection of information visualizations and graphic designs. The human perception studies in this thesis are run on information visualizations and graphic designs (Part III), which are also used for training a computational model of attention (Part IV), while the computer vision methods for detecting text and icons are developed with a focus on infographics (Part V).



additional importance annotation for a subset of 264 design variants from the *Design Improvement Results* (DIR) dataset [154] for evaluating our model of attention.

**Infographics:** We scraped and curated a dataset of 29K infographics from the *Visual.ly* website, a community platform for human-designed visual content. Chapter 11 introduces our *Visually29K* dataset that we used to train topic and category prediction, an icon proposal mechanism (Chapter 12), and for our multi-modal summary application (Chapter 13).

**Other:** We used static website images from the *FiWI* dataset [194] for evaluating our BubbleView interface (Chapter 7) and collected some academic poster images (from our colleagues, from the CVPR conference) for evaluating our ZoomMap interface (Chapter 8).

See Fig. 1.7 for some sample images from these datasets.

## ■ 1.3  Overview of thesis

This thesis has two complementary goals: (i) to develop data collection methods to measure how humans perceive multimodal documents (e.g., where people look, what they find important), and (ii) to develop computer vision tools to automatically parse

**Figure 1.7.** A sampling of the datasets used in this thesis: MIT300 [106] and OSIE [223] with natural images; CAT2000 [20] with images from different categories (including action, satellite, fractal images, as shown here); Visually29K (introduced in this thesis) with infographics; GDI [154] with graphic designs and DIR [154] with graphic design variants; MASSVIS [23] with visualizations; and FiWI [194] with webpages.



and make predictions about multimodal documents (e.g., what topics they are about).

Part I sets up the background for the rest of the thesis: discussing related past work on computational multimodal document understanding (Chapter 2), attention tracking and saliency (Chapter 3), and human perception of multimodal documents - with a focus on information visualizations (Chapter 4). Chapter 5 contains an introduction to the metrics that will be used for evaluation throughout the thesis.

Part II (Chapter 6) of this thesis includes a discussion of where people look in images and covers the state-of-the-art in computational models of visual attention. Models that predict a saliency value at each image pixel (as a measurement of the conspicuity of that pixel within the image) are called **saliency models**. Traditionally, saliency has been used to refer to bottom-up pop-out: a measure of whether regions of the image stand out from their surroundings/background, and whether they are likely to catch an observer's attention. Chapters 6 and 7 of this thesis argue for a re-branding of saliency as **importance**, showing that many regions on an image that an observer pays

attention to are semantically conspicuous, not necessarily visually conspicuous. For instance, human observers consistently look at a region of the image where they expect an object to be in the future (e.g., the landing position of a ball in mid-air) even if that region of the image has no distinctive visual features. Given a longer viewing interval, people spend time on regions of the image that they find interesting or important to understanding the scene. Extending this notion to non-natural images like graphic designs and information visualizations, people spend time on visual and textual elements that are important to understanding the underlying content. These attention patterns are not covered by the traditional definitions of saliency, nor can they be predicted by standard saliency models. In Part IV of the thesis, we build computational models of importance for graphic designs and information visualizations, the first fully-automatic and generalizable models of attention for these image types.

Data-driven models of attention have typically been trained on eye movements collected using an eye tracker. Recently, there has been a trend to train models on approximate attention data, collected as cursor data (instead of eye movements) in a crowdsourced online setting (e.g., [103]). This allows scaling up data collection to thousands of images, making the training of computational models with many parameters - like neural networks - more feasible. In a similar vein, for training the computational models of importance in Part IV, we designed a cursor-based interface called *Bubble-View* to collect importance data, which we introduce in Part III (Chapter 7) of the thesis. Part III of the thesis is about our data collection methods for crowdsourcing attention/importance. Our goal is to measure which regions of an image are particularly interesting to observers, when the observers are given a way of explicitly exploring or interacting with the images, by clicking or zooming (Chapters 7-8).

As discussed at the end of Chapter 7 (Part III) and in Chapter 10 (Part IV), there are many applications of importance prediction, ranging from design applications like automatic thumbnailing and retargeting, to image understanding tasks like captioning and visual question answering. Knowing what is most important to an observer in a document can help with summarization and retrieval of relevant content in search applications. Beyond predicting which regions of a design have highest importance, these applications require also an understanding of the underlying content itself (i.e., what is inside the design regions). In other words, to summarize or caption a design/document, we need to know what the underlying text elements are about and what the visual elements represent. The final part of this thesis, Part V, is about using text and icon parsing techniques to predict the topics that an infographic is about and to summarize the infographic. Finding that computer vision models trained on natural images did not generalize to graphics, we developed an icon proposal mechanism for infographics, trained on synthetic data (Chapter 12). We show that, taken together, text parsing, icon proposals, and icon classification can be used to densely annotate infographics, and this can serve as input to future applications for understanding designs (Chapter 13).

Overall, this thesis makes contributions to: human vision - by adding to the understanding of human attention on both natural and graphic images; user interfaces - by

introducing two novel data collection methods for capturing human attention; and to computer vision - by introducing approaches that can parse, annotate, and summarize non-natural images. Future directions along this line of research are to integrate the attention modeling and document understanding components developed in this thesis for captioning and summarization applications, for smarter designs tools (i.e., finding the most effective way to arrange textual and visual elements to communicate a message), and for automatically annotating text documents with images (e.g., for suggesting relevant visuals for articles and translating textbooks into interactive presentation slides) as in Fig. 1.3.

## ■ 1.4 Who is this thesis for?

I wrote this thesis with future students in mind: either students individually studying related topics who would like some additional background or inspiration, or students who may work with me on extending and building on the directions in this thesis. For this purpose I have written up future directions at the end of each chapter. The interdisciplinary nature of this thesis should cater to researchers in computer vision, human vision, and human-computer interaction; researchers in natural language processing, machine learning, and information visualization may gain additional ideas for their work. These research topics at the interface of multiple fields have a lot of unexplored potential and exciting applications, and it would be great to have a bigger community build on them.

## ■ 1.5  Papers included in this thesis

This section details the published papers that correspond to chapters in this thesis, with notes of how the paper was modified into the chapter.

**Chapter 5:** includes excerpts from the paper: Bylinskii, Z., Judd, T., Oliva, A., Torralba, A., Fredo, D. "What do different evaluation metrics tell us about saliency models?" [TPAMI 2018]. Specifically, the chapter consists of subsets of Sec. 3.2, Sec. 4, and a summary of Sec. 5.2, 5.4, and 7 from the paper. A section on the "Basics of human eye movements" was added to this chapter.

**Chapter 6:** is based in full on the paper: Bylinskii, Z., Recasens, A., Borji, A., Oliva, A., Torralba, A., Durand, F. "Where should saliency models look next?" [ECCV 2016] with some content added in from the Supplemental Material.

**Chapter 7:** is based on Kim, N.W.*, Bylinskii, Z.*, Borkin, M.A., Gajos, K., Oliva, A., Durand, F., Pfister, H. "BubbleView: an interface for crowdsourcing image importance maps and tracking visual attention" [TOCHI 2017]. Specifically, the chapter consists of Sec. 1, 3.2, 3.3, 5, 6, 7 and 8 from the paper. Sec. 7.5 has been added to the thesis after paper publication.

**Chapter 8:** is related to work that is in submission: Bylinskii, Z., Tancik, M., Newman, A., Zhong, K., Madan, S., Oliva, A., Durand, F. "ZoomMap: Using Zoom to Capture User Areas of Interest on Images".

**Chapters 9 and 10:** are based on Bylinskii, Z., Kim, N.W., O'Donovan, P., Alsheikh, S., Madan, S., Pfister, H., Durand, F., Russell, B., Hertzmann, A. "Learning Visual Importance for Graphic Designs and Data Visualizations" [UIST 2017]. Discussion of metrics has been omitted to avoid overlap with Chapter 5. Extra details have been added in from the Supplemental Material of the paper. Sec. 10.3 in this chapter is new.

**Chapters 11 and 13:** includes content from Bylinskii, Z.*, Alsheikh, S.*, Madan, S.*, Recasens, A.*, Zhong, K., Pfister, H., Durand, F., Oliva, A. "Understanding Infographics through Textual and Visual Tag Prediction" [arXiv:1709.09215 2017] (specifically, parts of Sec. 3, 4, and 5 from the paper) with major updates to Chapter 13 based on work that is currently in submission (as in Chapter 12).

**Chapter 12:** is related to work that is in submission: Bylinskii, Z.*, Madan, S*, Tancik, M.*, Zhong, K., Recasens, A., Alsheikh, S., Pfister, H., Durand, F., Oliva, A. "Synthetically Trained Icon Proposals for Parsing and Summarizing Infographics".

Additionally, a section on "Future directions" has been added at the end of every chapter to discuss limitations and natural extensions of the work presented.

# Part I

# Background

# Chapter 2

# Related work on multimodal document understanding

**C**OMPUTATIONAL approaches to processing multimodal documents have historically been tackled by different research communities: text parsing by the document understanding community, visualizations by the infovis community, websites by the human-computer interaction community, sketches and illustrations by the graphics community, etc. There has also recently been an increase in the number of papers and workshops [56, 121] on graphic design datasets and problems at top-tier computer vision conferences. The work presented in this thesis is timely in that it lies at the interface of computer vision, human perception, and human-computer interfaces, with applications to information visualizations and graphic designs. In this chapter we discuss relevant papers on multimodal document understanding coming from these different research communities.

**Document analysis:** *ICDAR*, the International Conference on Document Analysis and Recognition, is a biannual conference on character and text recognition, document analysis and understanding. Relevant work includes classifying documents by type (e.g.,

**Figure 2.1.** Figures from the ICDAR Competition on Recognition of Documents with Complex Layouts [5] (http://primaresearch.org/RDCL2017/). The competition consists of segmenting pages of magazines and technical articles, classifying regions (text, image, table, chart, math, etc.), and extracting and recognizing the contained text.

**Figure 2.2.**    Figure from ReVision by Savva et al. [191] demonstrating the automatic application of converting a bitmap visualization into a redesigned editable chart, after parsing all the visualization components (text, values, axes, marks, etc.) and converting them into a structured table representation.



email, news article, presentation, scientific publication) [79], separating figures from text in articles [213] and more fine-grained region classification problems [5], where document regions are labeled as text, image, graphic, table, math, etc. (Fig. 2.1). Related are vision-based and DOM-based approaches that decompose a website into sub-parts for further analysis [37]. A separate class of approaches transcribe the text from document pages into characters (i.e., optical character recognition methods) [200]. Most document analysis and retrieval methods, however, stop short of processing the semantics of the document elements [144]. Document analysis can therefore be used as pre-processing for computer vision methods to understand and classify the content inside the images, and for natural language processing (NLP) techniques to make sense of the parsed text - for topic understanding, summarization, and question-answering applications.

**Figure parsing:** Computer vision methods have been applied to classify scholarly figures and graphs [170, 191], to transcribe textbook diagrams into structured tables for question answering [109, 193], and to parse the elements (axes, legends, plot lines, etc.) within graphs and charts for re-design applications [92, 167, 191, 196] (Fig. 2.2). Siegel et al. [196] show that their figure parsing results can be used to answer queries and summarize results (but their results are limited to 2D-graphical plots of precision-recall and ROC curves presented in technical papers). Wu et al. [221] use the XML representation of a line graph to predict the intended message of the graph, and cite as motivation the need to summarize the content of multimodal documents by taking into account the contained information graphics as well as the text. Other approaches have looked at converting hand-drawn sketches into computational diagrams, requiring an automatic parsing and classification of sketches into common diagrammatic components [4, 53], as an automatic design aid for paper figure generation.

**Machine learning for design:** Our work also relates to the general program of applying machine learning in the service of graphic design tools. *Webzeitgeist* [127] is a repository and platform of 100K webpages and 100M design elements and *RICO* [47] is a dataset of 9.7K mobile apps covering 72K unique UI screens, both datasets collected to enable large-scale statistical analysis of design patterns and design-driven search and

**Figure 2.3.** Figure from SmartNails by Berkner et al. [12] demonstrating the proposed thumbnails (bottom) automatically-computed from the document images (top) by cropping, scaling, and re-arranging image and text elements.



machine learning. Ritchie et al. [183] introduce another style-based exploration tool to mine a design gallery for suggestions. Saleh et al. [188] learn a style similarity metric for infographics for image retrieval tasks using a newly collected dataset of 19K infographics and crowdsourced similarity data. With a similar application of style-based retrieval from a repository, Garces et al. introduce a similarity metric for illustrations [61].

**Related work in computer vision:** Computer vision has traditionally focused on understanding natural images and scenes. However, there is a growing interest in graphic designs, which motivates a new set of research questions and technical challenges. Zitnick et al. [232, 233] introduced abstract scenes to study higher-level image semantics (relationships between objects, storylines, etc.). Wilber et al. [218] presented an 'Artistic Media Dataset' to explore the representation gap between objects in photographs versus in artistic media. Iyyer et al. [98] built a 'COMICS' dataset and made predictions about actions and characters using extracted visual and textual elements from comic panels. Hussain et al. [95] presented a dataset of advertisements and described the challenges of parsing symbolism, memes, humor, and physical properties from images. Other relevant work to the problem of abstract image recognition and representation is work on sketches [76, 189]. To the best of our knowledge, there is no computer vision work on automated understanding of infographics.

**Retargeting and summarization of documents:** Previous work has explored the idea of summarizing multimodal content in thumbnail form [12, 203, 220]. Berkner et al. [12] crop, scale, and re-arrange visual and textual elements into dense summaries of the content called 'SmartNails' (Fig. 2.3). Importance of a visual region is measured using the bit-rate of the compressed JPEG representation. Importance of a text region is measured as a product of the font size and position. A greedy algorithm is used to arrange the important visual and textual regions into the final thumbnail. Strobelt et al. [203] combine images and key terms from an entire document to create 'Document Cards'. Key terms are selected based on a frequency analysis, and images are considered important if key terms are found in their captions. A packing algorithm

**Figure 2.4.** Figure from DesignScape by O'Donovan et al. [155] demonstrating a system which aids the design process by making interactive layout suggestions (in the panels to the left and right of the central canvas) as a user moves around design elements on a canvas.



is then used for arranging the elements in the final thumbnails. Woodruff et al. [220] utilize user feedback to create thumbnails of webpages, by modifying and amplifying HTML elements. The effectiveness of the proposed thumbnails is demonstrated using a search application. Related also is work on automated retargeting where the elements of a graphic design are re-arranged for a different form factor. O'Donovan et al. takes a machine learning approaching in combination with design heuristics to adjust the layout of graphic designs, with and without user input [154, 155] (Fig. 2.4). Kumar et al. [126] retarget webpages designed for the desktop to mobile devices using the underlying content and style files.

Chapter 13 presents a purely machine learning approach to the multi-modal summarization of infographics (see Fig. 1.4 for an example). The text in infographics is first automatically extracted using out-of-the-box OCR techniques [68], and key text topics are predicted using a trained bag-of-words to topic model (a neural network). The visual regions (i.e., icons) in infographics are extracted using a trained icon proposal mechanism (Chapter 12). The final visual summary is composed by choosing the top 1-3 predicted text topics, and for each text topic, the most visually relevant icon is selected (automatically, using a neural net classifier). Compared to prior work on summarization, this approach is maximally data-driven and does not rely on heuristics or design decisions.

In Chapter 10, we take a perception based approach (with computer vision models) to image processing and summarization. Visualizations and graphic designs are treated as bitmap images. Unlike a lot of the related work, no vector input or other structured information (e.g., DOM structure, style files, annotated regions, user-selected elements) are required. Attention maps are predicted automatically on a per-pixel level, and the resulting applications use the attention maps directly to thumbnail and retarget the underlying designs. This is inspired by saliency-based retargeting of natural images [6, 186]. The work in this thesis presents the first fully-automatic and generalizable (to different visualization and graphic design styles) model of saliency, that is used for automatic thumbnailing and retargeting. The next chapter discusses work on attention and saliency modeling that is also relevant to this approach.

# Chapter 3

# Related work on attention tracking and saliency

**T**HE eyes are a window onto the mind. A significant amount of research has demonstrated connections between eye movements and various cognitive tasks: the eyes can provide important clues about how visual perception proceeds as a human looks at images [83, 85, 108, 122, 152]. This area of research is so established and diverse that we refer the reader to some representative papers reporting on the utility of eye movements for studying human perception and cognition in the context of user interfaces [11, 27, 49, 66, 70, 99, 169, 180], web search [43, 67], web browsing [42, 105, 156], problem solving [71], reading [172], advertisements [174], and visualizations [24, 34, 90, 115, 168]. These papers show that aside from providing information about how human perception proceeds, eye movements can also provide insights about the effectiveness of different visual content, or the usability of interfaces.

Motivated by the implications of knowing where people look, many researchers have sought ways to efficiently collect eye movements at-scale, often as an alternative to more standard eye tracking procedures. In this chapter we cover work in this area most relevant to our own user interfaces for capturing attention - Part III of the thesis. As background to Parts II and IV of the thesis, this chapter also includes some brief discussion on saliency models, which are computational models of visual attention.

## ■ 3.1 Cursor-based attention tracking

There has been a significant effort to find cheap, nonintrusive, and more scalable alternatives to collect human attentional data. Cursor-based techniques are a particularly suitable alternative for scaling to large web-based studies.

The **moving-window** approach is a popular cursor-based technique in which a limited amount of information is visible through a variable size window continuously following a cursor position [146, 173]. Inspired by the moving-window model, Jansen et al. [101] developed a computer program called Restricted Focus Viewer (RFV) that takes an image, blurs it, and reveals only a restricted block of the image, allowing a user to move the region using a mouse [10, 14, 101, 205] (Fig. 3.1). Commercial software for tracking user attention has also built on the same idea. For instance, Attensee

(www.attensee.com) is a commercial solution based on the idea of Flashlight [192], an open-source research tool (github.com/michaelschulte/flashlight). The mouse-contingent methodology has been employed to investigate cognitive behaviors of users in diverse contexts such as diagrammatic reasoning and program debugging, and to study the usability of web sites [9, 101, 205].

**Figure 3.1.** Example figure from the Restricted Focus Viewer by Jansen et al. [101], a moving-window interface that allows participants to explore a blurred diagram through a focus window by moving the computer cursor (cursor-based attention tracking).

Recent studies have made further improvements. SALICON [103] implemented moving-window, multi-resolution blur on images to attempt to simulate the fall-off in acuity of peripheral vision. On the other hand, Lagun and Agichtein [132] directly preprocessed web search results to show one result and blur the other results based on a user's viewport; however, this method is not intended to approximate the human fovea as it shows an entire DOM element at a time. All these recent studies were conducted online with hundreds to thousands of participants, proving the scalability of their methods.

There is also a rich history of work in the space of gaze-contingent multiresolutional displays, where the moving-window approach is guided by gaze. We refer the reader to a review by Reingold et al. [178]. These approaches complement, rather than replace, standard eye-tracking techniques, and have different motivations: bandwidth and processing savings. However, this line of work contains an investigation of multiresolutional blur to approximate the peripheral visual system, which is relevant to the analyses in our studies (Chapter 7). Whether cursor-based or gaze-based, a moving-window approach slows down visual exploration patterns relative to natural viewing and can be used to discover the most important or relevant image regions.

Aside from the moving-window model for image exploration, other works also investigated the relationship between cursor movements and gaze positions, mostly focusing on web browsing [39] and search tasks [74, 87, 88, 184]. Chen et al. [39] found a high correlation between cursor and gaze locations. Rodden et al. [184] found that cursor and gaze are better aligned along the vertical dimension, while Guo and Agichtein [74] also found a similar result in their study of predicting eye-mouse coordination. Huang

et al. [88] found that people's cursors lag behind their gazes and there are individual differences in the distance between the cursor and gaze positions.

## ■ 3.2  Appearance-based gaze tracking

Another line of work has been devoted to non-intrusive, appearance-based gaze estimation, where images of the eyes are post-processed using computer vision techniques to determine gaze location. This type of gaze estimation often involves collecting a training dataset with a standard eye-tracker, training a computer vision model to map eye images to gaze coordinates, and using this model at test-time to directly infer gaze positions from a video stream of the eyes (e.g., captured via a webcam). At test time, these approaches do not require specialized eye tracking hardware (i.e., high quality special cameras, infrared sensors, and head mounting devices).

Early gaze tracking models were mostly based on relatively small training datasets collected through lab studies. For example, Baluja and Pomerleau [7] collected 2000 images of the eyes for four postures by instructing a participant to visually track a moving cursor and built a neural network model to estimate gaze locations. Recent methods attempt to build gaze tracking models on large datasets to improve accuracy as well as to work in real-world settings. Funes Mora et al. [59] constructed a database to enable comparison across different gaze tracking algorithms for variations including head poses, individual differences, and ambient and sensing conditions. Zhang et al. [229] developed an appearance-based gaze estimation method using multimodal convolutional neural networks. Their model was trained on a hundred thousand images from 15 laptop users for several months using built-in cameras in laptops, accounting for realistic variability in illumination and appearance. Huang et al. [89] similarly built a large gaze dataset and a gaze tracking algorithm for tablet users. While the two studies are still limited to datasets collected through labs, other works leverage online crowdsourcing to further extend the scale of gaze datasets. Xu et al. [224] developed a webcam-based eye tracking game running in a browser on a remote computer. Their crowdsourced experiments could collect gaze data cheaper and faster than lab studies. Papoutsaki et al. [160] also designed a similar webcam-based eye tracking system. Krafka et al. [123] collected eye tracking on over 2.5M frames using a mobile application and online, and developed a gaze prediction algorithm based on convolutional neural networks, while achieving state-of-the-art results (Fig. 3.2).

All of the above approaches have yet to reach the level of tracking accuracy and robustness possible with dedicated eye tracking hardware. These approaches also depend on either some initial calibration or have constraints on a participants' set-up: network connection, camera quality, and restricted range of face location relative to screen. As a result, we have not yet seen widespread adoption of appearance-based gaze tracking. Additionally, the camera-based gaze tracking approaches have the downside of requiring the capture of participants' face images throughout the study, which comes with privacy concerns [136].

**Figure 3.2.**   Example figure from iTracker by Krafka et al. [123], a convolutional neural network for predicting eye gaze coordinates from face and eye images captured using mobile cameras (appearance-based gaze tracking).   This network was trained on 2.5M frames of 1450 crowdsourced participants performing a calibration task.



## ■ 3.3  Saliency models and datasets

In addition to alternative techniques for eye tracking which require human participants, significant progress has been made building computational saliency models to predict eye fixations. Many saliency models are motivated by psychological and neurobiological theories, and make use of both low-level image features (e.g., intensity, color, and orientation) and high-level semantic features (e.g., scenes, objects, and tasks) to approximate the human visual system [17, 58]. The performance of these models is usually evaluated against ground-truth eye fixations [31, 36, 107].

**Traditional models of saliency:** Computational modeling of bottom-up attention dates back to the seminal works by Treisman and Gelade [212] (Feature Integration Theory), the computational architecture by Koch and Ullman [117] and the bottom-up model of Itti et al. [96, 97] (Fig. 3.3). Parkhurst and Neibur were the first to measure saliency models against human eye fixations in free-viewing tasks [161]. Followed by this work and the Attention for Information Maximization model of Bruce and Tsotsos [26], a cascade of saliency models emerged, establishing saliency as a subarea in computer vision. Large datasets of human eye movements were constructed to provide training data, object detectors and scene context were added to models, and learning approaches gained traction for discovering the best feature combinations [19, 106, 111, 223, 230]. Please refer to [17, 18] for recent reviews of saliency models.

**Neural network models of saliency:** One of the first attempts to leverage deep learning for saliency prediction was Vig et al. [215], using convnet layers as feature maps to classify fixated local regions. Kümmerer et al. [129] introduced the model DeepGaze, built on top of the AlexNet image classification network [124]. Similarly, Liu et al. [139] proposed the Multiresolution-CNN model in which three convnets, each on a different image scale, are combined to obtain the saliency map. In the SALICON model [103], CNNs are applied at two different image scales: fine and coarse.

Traditionally, saliency models have been trained directly on fixation data collected from eye tracking experiments [106, 111]. However, deep neural network models require large quantities of data, larger than what is practical to collect using conventional eye tracking techniques. To overcome this challenge, a large-scale crowd-sourced dataset of mouse movements on natural images was recently released for simulating the nat-

**Figure 3.3.** Figure from Itti and Koch [96] showing a traditional saliency architecture, where different feature maps, capturing orientation, intensity, and colors, are extracted from an image at multiple scales, center-surround differences are computed to obtain locally salient regions, and the resulting activations are combined into a final saliency map.

ural viewing behavior and subsequently training computational saliency models. This dataset, dubbed SALICON, was collected using a moving-window methodology [103]. Since then, many neural network models of saliency trained on this data [103, 125, 157] have achieved state-of-the-art performances on standard saliency benchmarks [31].

For instance, DeepFix [125] is a fully convolutional neural network (convnet) built on top of the VGG network [198] and trained on the SALICON dataset to predict pixel-wise saliency values in an end-to-end manner. DeepFix has additionally been fine-tuned on MIT1003 [106] and CAT2000 [20]. Pan et al. [157] trained two architectures on SALICON in an end-to-end manner: a shallow convnet trained from scratch, and a deeper one whose first three layers were adapted from the VGG network (Sal-Net). Tavakoli et al. [209] have recently shown that saliency models trained on mouse movements can generalize well to predicting eye fixations. Other saliency models based on deep learning have been proposed for salient region detection [133, 135, 217, 231]. While deep learning models have shown impressive performance for saliency prediction, a finer-grained analysis shows that they continue to miss key elements in images. In Chapter 6 of this thesis we discuss the advances in saliency models, remaining gaps to human performance, and new directions forward to more closely approximate human attention.

## ■ 3.4  Saliency models for non-natural images

While most saliency models are focused on predicting eye fixations on natural scenes, there are relatively few studies that have looked at other image types including web-pages, graphic designs, and information visualizations. These images are different from natural images in that they usually contain rich semantic data (e.g., texts, charts, and logos) or different viewing patterns such as top-left bias [29] and banner blindness [73]. Designers and researchers have long studied eye movements as a clue to understanding

**Figure 3.4.**  Figure from Pang et al. [159] demonstrating a web design application where, given a webpage as input (left panel), a designer specifies a desired trajectory over page components (middle panel) and an automatic system modifies the position, style, and color of the components (right panel) such that the predicted eye gaze patterns match the trajectory as closely as possible.



Input web design                    Input path                    Output web design

the perception of user interfaces [50, 99]. There have also been several recent studies of eye movements and the perception of designs [24, 82]. However, few researchers have attempted to automatically predict attention in graphic designs.

The DesignEye system [185] uses hand-crafted saliency methods, demonstrating that saliency methods can provide valuable feedback in the context of a design application. O'Donovan et al. [154] gather crowdsourced importance annotations, where participants are asked to mask out the most important design regions. They train a predictor from these annotations. Haass et al. [77] test three natural image saliency models on the MASSVIS data visualizations [23], concluding that, across most saliency metrics, these models perform significantly worse on visualizations than on natural images. Several models also exist for web page saliency. However, most methods use programmatic elements (e.g., the DOM) as input to saliency estimation [29, 220]. Pang et al. predict the order in which people will look at components on a webpage [159] by making use of the DOM and manual segmentations (Fig. 3.4). Other works use the web page image itself as input to predict saliency [194, 202]. For instance, Shen and Zhao [194] developed a webpage saliency model based on the FiWI dataset, and then improved the model with high-level semantic features (e.g., positional bias and object detectors) [195]. Xu et al. [225] presented a computational model for predicting visual attention in user interfaces with user interactions.

Inspired by the work of O'Donovan et al. [154], in Chapter 9 we present a fully automated model of importance for both graphic designs, but also information visualizations. While trained on different datasets, the model architecture is the same for both image types. Unlike O'Donovan et al., however, our method does not require knowledge of the location of design elements (i.e., a vector representation) to run on a new design. The input to our model is the bitmap image of the design alone. Moreover, our model is based on a neural network architecture which gives it predictive power over previous approaches (to predicting saliency on non-natural images) that are based on lower-level features and older saliency models. The neural network architecture also runs at 0.1 seconds per image, making it 100 times faster than O'Donovan et al.'s model, and practical for interactive design applications (Chapter 10, Sec. 10.4).

# Chapter 4

# Related work on perception of information visualizations

**U**NDERSTANDING how people perceive information visualizations - what catches their attention, what they spend time looking at, and what they ultimately remember - can both help us design computational systems to process visualizations, and to better understand what makes a visualization effective for design applications.

Many important works in the visualization community have studied how different visualization types are perceived, and the effect of different data types and tasks [40, 120, 166]. The effect of "visual embellishments" on the memorability and comprehensibility of visualizations is also an active area of research [8, 13, 16, 23, 94, 134, 145, 199, 214]. The effect and role of specific visual elements have also been investigated within the context of specific visualization types, e.g., attributes of node link diagrams [3, 62], specific visual elements such as pictographs [81], visual distortions [158], and more broadly [80].

Apart from memory, eyetracking has also been used to study how a person views and visually explores a visualization [15]. Previous studies have focused on eye movements on specific visualization types such as graphs [90, 91, 119, 168], tree diagrams [28], and parallel coordinates [197], for comparing multiple types of visualizations [65], and for evaluating cognitive processes in visualization interactions [115]. There has also been research in the area of understanding different types of tasks and visual search strategies for visualizations through the analysis of eyetracking fixation patterns as well as insights into cognitive processes [168, 171].

In our own past work, we investigated the human perception of information visualizations, including more traditional graph types, as well as scientific visualizations and infographics from the MASSVIS dataset [23, 24]. We ran studies online and in the lab to measure human attention and memory. We collected eye movements using an in-lab eyetracking set-up, memorability scores using both crowdsourced web experiments and our in-lab set-up, and text descriptions.

Running large-scale memorability studies with hundreds of visualizations and dozens of participants, we learned which features of visualizations increase their memorability - i.e. whether they'd be remembered if shown to participants at a later time point [23].

**Figure 4.1.** We measured which visualizations people remembered best, as well as which visualizations were best recalled from memory (measured by the quality of participant-generated descriptions). Memorable visualizations tended to be visually distinct, dense, and colorful. Effective visualizations made good use of visual elements (pictograms) and had informative titles.



We found that visualizations with pictograms or recognizable objects were more memorable, as were visualizations that were visually distinct, colorful and visually dense (Fig. 4.1). Visually distinct visualizations are those that are represented in unique ways - e.g., a bar graph where the bars are composed of stacked coins. That visualizations with many visual elements are more memorable is not surprising. What is surprising is that people not only better recognized these visualizations, but could better retrieve details about the visualizations from their memory [24]. In other words, people could better recall facts and information encoded in a visualization if a visualization was more visually memorable. Visual elements like pictographs act like hooks into memory, helping to retrieve the information encoded in the visual associations made. By visually representing certain concepts in visualizations, the effectiveness of those visualizations (measured as memory for retrieved details) improves. This should inform both how visualizations should be designed, and the importance that computational algorithms place on detecting visual elements in understanding visualizations (Part V of thesis).

By analyzing eye movements [24], we found that the visualization elements that people spend the most time visually exploring are the text elements, and especially the titles of visualizations (Fig. 4.2). Those elements influenced what participants subsequently recalled. In other words, knowing where people look can provide us with important clues about what they encode into, and can later recall from, memory. We can measure where people look when they're exploring a visualization for the first time, and where they look on a visualization to remember if they have seen it previously (Fig. 4.3). This helps identify which parts of a visualization trigger the memory. A conclusion of these analyses is that the measurement of eye movements can provide a way to investigate human memory, comprehension, and can also be used to analyze the effectiveness of different designs (which serves as motivation for Parts III-IV of the thesis).

**Figure 4.2.** Relative importance of different visualization elements, measured as the density of participant eye fixations on those elements, averaged over hundreds of visualizations. Participants spent the most time visually exploring the text elements on visualizations, particularly the titles.

**Figure 4.3.** Examples of the most and least recognizable visualizations. TOP: Eye-tracking fixation heat maps (i.e., average of all participants' fixation locations) when participants viewed the visualizations for the first time. The fixation patterns demonstrate visual exploration of the visualization. BOTTOM: Eye-tracking fixation heat maps when participants viewed visualizations for a second time and were asked if they remember them. The most recognizable visualizations all have a single focus in the center indicating quick recognition of the visualization, whereas the least recognizable visualizations have fixation patterns indicative of visual exploration (e.g., title, text, etc.) for recognition.

# Chapter 5

# Metrics for evaluating saliency models

*This chapter consists of excerpts from: Bylinskii, Z., Judd, T., Oliva, A., Torralba, A., Fredo, D. "What do different evaluation metrics tell us about saliency models?" [TPAMI 2018]*

**I**N this chapter, we introduce the saliency metrics that will be used for evaluation throughout the rest of this thesis. They will be used to evaluate saliency models (Chapter 6), to compare attention maps captured using alternative user interfaces to eye fixations (Chapters 7 and 8), and to evaluate our computational models at the ability to predict importance on information visualizations and graphic designs (Chapter 9).

Saliency metrics are functions that take two inputs representing eye fixations (ground truth and predicted) and then output a number assessing the similarity or dissimilarity between them. Given a set of ground truth eye fixations, such metrics can be used to define scoring functions, which take a saliency map prediction as input and return a number assessing the accuracy of the prediction.

We consider 8 popular saliency evaluation metrics in their most common variants. Some metrics have been designed specifically for saliency evaluation (shuffled AUC, Information Gain, and Normalized Scanpath Saliency), while others have been adapted from signal detection (variants of AUC), image matching and retrieval (Similarity, Earth Mover's Distance), information theory (KL-divergence), and statistics (Pearson's Correlation Coefficient). Because of their original intended applications, these metrics expect different input formats: KL-divergence and Information Gain expect valid probability distributions as input, Similarity and Earth Mover's Distance can operate on unnormalized densities and histograms, while Pearson's Correlation Coefficient (CC) treats its inputs as random variables. Following Riche et al. [181], we categorize metrics as **location-based** or **distribution-based**, depending on whether the ground-truth is represented as discrete fixation locations or a continuous fixation map, accordingly. Visualizations of all these metrics to add intuition to their computation can be found in Bylinskii et al. [36].

## ■ 5.1 Basics of human eye movements

The human eye consists of light receptor cells that are differently distributed throughout the eye. The clearest and most detailed vision is in the central, **foveal area**, of the visual field, and blurrier vision is in the larger part of the visual field, which is called the **peripheral area**. The foveal area captures about 1-2 degrees of visual angle which constitutes less than 8% of the visual field, but makes up 50% of the visual information sent to the brain [211]. When we move our eyes, we place the foveal region of the eye on different regions of the visual field, bringing them into focus.

    **Visual angles** are units for measuring the projection of the visual field, as images, on our retina. For a given experimental viewing setup, visual angles can be computed by taking into account the distance to the screen, size and resolution of the image on the screen[1]. The error of professional-grade eye trackers (e.g., EyeLink) is also measured in degrees of visual angle, and is commonly less than 1 degree.

    The pauses in eye movements are called **fixations**, and the transitions between successive fixations are called **saccades**. In all the work described in this thesis, we focus on fixations, since they give us the points of interest that the eye has stopped on to bring them into focus. The temporal sequence of fixations, fixation duration, saccade length, and other features of eye movements carry a lot of additional information about human perception [34, 85, 99, 108] but are beyond the scope of the present work. We concentrated on the location of fixations, which are most straightforward to analyze [27, 99, 211] and to model computationally [36].

## ■ 5.2 Using collected eye fixations as ground truth for evaluation

Ground truth eye fixations can be processed and formatted in a number of ways for saliency evaluation. There is a fundamental ambiguity in the correct representation for the fixation data, and different representational choices rely on different assumptions. One format is to use the original fixation locations. Alternatively, the discrete fixations can be converted into a continuous distribution, a **fixation map**, by smoothing. We follow common practice and blur each fixation location using a Gaussian with sigma equal to one degree of visual angle [147]. Throughout this chapter, we denote the map of fixation locations as $Q^B$ and the continuous fixation map (distribution) as $Q^D$.

    Smoothing the fixation locations into a continuous map acts as regularization. It allows for uncertainty in the ground truth measurements to be incorporated: error in the eye-tracking as well as uncertainty of what an observer sees when looking at a particular location on the screen. Any splitting of observer fixations in two sets will never lead to perfect overlap (due to the discrete nature of the data), and smoothing provides additional robustness for evaluation. In the case of few observers, smoothing the fixation locations helps to extrapolate the existing data. The fixation locations can be viewed as a discrete sample from some ground truth distribution that the fixation map attempts to

---

[1] https://github.com/cvzoya/saliency/tree/master/computeVisualAngle

approximate. Similarly, the fixation map can be viewed as an extrapolation of discrete fixation data to the case of infinite observers.

## ■ 5.3 Location-based metrics

## ■ 5.3.1 Area under ROC Curve (AUC)

**Evaluating saliency as a classifier of fixations:** Given the goal of predicting the fixation locations on an image, a saliency map can be interpreted as a classifier of which pixels are fixated or not. This suggests a detection metric for measuring saliency map performance. In signal detection theory, the Receiver Operating Characteristic (ROC) measures the tradeoff between true and false positives at various discrimination thresholds [57, 72]. The Area under the ROC curve, referred to as AUC, is the most widely used metric for evaluating saliency maps. The saliency map is treated as a binary classifier of fixations at various threshold values (level sets), and an ROC curve is swept out by measuring the true and false positive rates under each binary classifier (level set). Different AUC implementations differ in how true and false positives are calculated. In the implementation we use from  Judd et al. [106], the true positive rate is the ratio of true positives to the total number of fixations, where true positives are saliency map values above threshold at *fixated pixels*. This is equivalent to the ratio of fixations falling within the level set to the total fixations. The false positive rate is the ratio of false positives to the total number of saliency map pixels at a given threshold, where false positives are saliency map values above threshold at *unfixated pixels*. This is equivalent to the number of pixels in each level set, minus the pixels already accounted for by fixations. Visualizations of AUC can be found in Bylinskii et al. [36].

**Penalizing models for center bias:** The natural distribution of fixations on an image tends to include a higher density near the center of an image [206]. As a result, a model that incorporates a center bias into its predictions will be able to account for at least part of the fixations on an image, independent of image content. In a center-biased dataset, a center prior baseline will achieve a high AUC score. The shuffled AUC metric, **sAUC** [21, 52, 206, 207, 228] samples negatives from fixation locations from other images, instead of uniformly at random. This has the effect of sampling negatives predominantly from the image center because averaging fixations over many images results in the natural emergence of a central Gaussian distribution [206, 219]. A model that only predicts the center achieves an sAUC score of 0.5 because at all thresholds this model captures as many fixations on the target image as on other images (true positive and false positive rates are equal). A model that incorporates a center bias into its predictions is putting density in the center at the expense of other image regions. Such a model will score worse according to sAUC compared to a model that makes off-center predictions, because sAUC will effectively discount the central predictions. In other words, sAUC is not invariant to whether the center bias is modeled: it specifically penalizes models that include the center bias.

**Invariance to monotonic transformations:** AUC metrics measure only the relative (i.e., ordered) saliency map values at ground truth fixation locations. In other words, the AUC metrics are ambivalent to monotonic transformations. AUC is computed by varying the threshold of the saliency map and comparing true and false positives. Lower thresholds correspond to measuring the coverage similarity between distributions, while higher thresholds correspond to measuring the similarity between the peaks of the two maps [54]. Due to how the ROC curve is computed, the AUC score for a saliency map is mostly driven by higher thresholds: i.e., the number of ground truth fixations captured by the peaks of the saliency map, or the first few level sets. Models that place high-valued predictions at fixated locations receive high scores, while low-valued predictions at non-fixated locations are mostly ignored.

### ■ 5.3.2 Normalized Scanpath Saliency (NSS)

**Measuring the normalized saliency at fixations:** The Normalized Scanpath Saliency, **NSS**, was introduced to the saliency community as a simple correspondence measure between saliency maps and ground truth, computed as the average normalized saliency at fixated locations [163]. Unlike in AUC, the absolute saliency values are part of the normalization calculation. NSS is sensitive to false positives, relative differences in saliency across the image, and general monotonic transformations. However, because the mean saliency value is subtracted during computation, NSS is invariant to linear transformations like contrast offsets. Given a saliency map $P$ and a binary map of fixation locations $Q^B$:

$$NSS(P, Q^B) = \frac{1}{N} \sum_i \overline{P_i} \times Q_i^B$$
$$\text{where } N = \sum_i Q_i^B \text{ and } \overline{P} = \frac{P - \mu(P)}{\sigma(P)}$$

(5.1)

where $i$ indexes the $i^{th}$ pixel, and $N$ is the total number of fixated pixels. Chance is at 0, positive NSS indicates correspondence between maps above chance, and negative NSS indicates anti-correspondence. For instance, a unity score corresponds to fixations falling on portions of the saliency map with a saliency value one standard deviation above average.

Recall that a saliency model with high-valued predictions at fixated locations would receive a high AUC score even in the presence of many low-valued false positives. However, false positives lower the normalized saliency value at each fixation location, thus reducing the overall NSS score.

### ■ 5.3.3 Information Gain (IG)

**Evaluating information gain over a baseline:** Information Gain, **IG**, was recently introduced by Kümmerer et al. [128, 130] as an information theoretic metric that mea-

sures saliency model performance beyond systematic bias (e.g., a center prior baseline). Given a binary map of fixations $Q^B$, a saliency map $P$, and a baseline map $B$:

$$IG(P, Q^B) = \frac{1}{N} \sum_i Q_i^B [\log_2(\epsilon + P_i) - \log_2(\epsilon + B_i)] \qquad (5.2)$$

where $i$ indexes the $i^{th}$ pixel, $N$ is the total number of fixated pixels, $\epsilon$ is for regularization. This metric measures the average information gain of the saliency map over the center prior baseline at fixated locations (at $Q^B = 1$), in bits per fixation.

IG assumes that the input saliency maps are probabilistic, properly regularized and optimized to include a center prior [128, 130]. A score above zero indicates the saliency map predicts the fixated locations better than the center prior baseline. This score measures how much image-specific saliency is predicted beyond image-independent dataset biases, which in turn requires careful modeling of these biases. We refer the reader to [130] for detailed discussions of the IG metric.

## ■ 5.4  Distribution-based metrics

The (location-based) metrics described so far measure the accuracy of saliency models at predicting discrete fixation locations. If the ground truth fixation locations are interpreted as a sample from some underlying probability distribution, then another approach is to predict the distribution directly instead of the fixation locations. Although we can not directly observe this ground truth distribution, it is often approximated by Gaussian blurring the fixation locations into a fixation map. Next we discuss a set of metrics that measure the accuracy of saliency models at approximating the continuous fixation map.

### ■ 5.4.1  Similarity (SIM)

**Measuring the intersection between distributions:** The similarity metric, **SIM** (also referred to as *histogram intersection*), measures the similarity between two distributions, viewed as histograms. First introduced as a metric for color-based and content-based image matching [187, 204], it has gained popularity in the saliency community as a simple comparison between pairs of saliency maps. SIM is computed as the sum of the minimum values at each pixel, after normalizing the input maps. Given a saliency map $P$ and a continuous fixation map $Q^D$:

$$SIM(P, Q^D) = \sum_i \min(P_i, Q_i^D)$$
$$\text{where } \sum_i P_i = \sum_i Q_i^D = 1 \qquad (5.3)$$

iterating over discrete pixel locations $i$. A SIM of one indicates the distributions are the same, while a SIM of zero indicates no overlap. SIM is very sensitive to missing values, and penalizes predictions that fail to account for all of the ground truth density.

The SIM metric is good for evaluating partial matches, where a subset of the saliency map accounts for the ground truth fixation map. As a side-effect, false positives tend to be penalized less than false negatives. For other applications, a metric that treats false positives and false negatives symmetrically, such as CC or NSS, may be preferred.

### ■ 5.4.2 Pearson's Correlation Coefficient (CC)

**Evaluating the linear relationship between distributions:** The Pearson's Correlation Coefficient, **CC**, also called *linear correlation coefficient* is a statistical method used in the sciences to measure how correlated or dependent two variables are. Interpreting saliency and fixation maps, $P$ and $Q^D$, as random variables, CC measures the linear relationship between them [148]:

$$CC(P, Q^D) = \frac{\sigma(P, Q^D)}{\sigma(P) \times \sigma(Q^D)} \tag{5.4}$$

where $\sigma(P, Q^D)$ is the covariance of $P$ and $Q^D$. CC is symmetric and penalizes false positives and negatives equally. It is invariant to linear (but not arbitrary monotonic) transformations. High positive CC values occur at locations where both the saliency map and ground truth fixation map have values of similar magnitudes. Due to its symmetric computation, CC can not distinguish whether differences between maps are due to false positives or false negatives. Other metrics may be preferable if this kind of analysis is of interest.

### ■ 5.4.3 Kullback-Leibler divergence (KL)

**Evaluating saliency with a probabilistic interpretation:** Kullback-Leibler (**KL**) is a broadly-used information theoretic measure of the difference between two probability distributions. In the saliency literature, depending on how the saliency predictions and ground truth fixations are interpreted as distributions, different KL computations are possible. We discuss a few alternative varieties in the appendix. To avoid future confusion about the KL implementation used, we can refer to this variant as **KL-Judd** similarly to how the AUC variant traditionally used on the MIT Benchmark is denoted AUC-Judd. Analogous to our other distribution-based metrics, KL-Judd takes as input a saliency map $P$ and a ground truth fixation map $Q^D$, and evaluates the loss of information when $P$ is used to approximate $Q^D$:

$$KL(P, Q^D) = \sum_i Q_i^D \log \left( \epsilon + \frac{Q_i^D}{\epsilon + P_i} \right) \tag{5.5}$$

where $\epsilon$ is a regularization constant[2]. KL-Judd is an asymmetric dissimilarity metric, with a lower score indicating a better approximation of the ground truth by the saliency map.

---

[2]The relative magnitude of $\epsilon$ will affect the regularization of the saliency maps and how much zero-valued predictions are penalized. The MIT Saliency Benchmark uses MATLAB's built-in *eps* = 2.2204e-16.

## ■ 5.4.4 Earth Mover's Distance (EMD)

**Incorporating spatial distance into evaluation:** All the metrics discussed so far have no notion of how spatially far away the prediction is from the ground truth. Accordingly, any map that has no pixel overlap with the ground truth will receive the same score of zero, regardless of how predictions are distributed. Incorporating a measure of spatial distance can broaden comparisons, and allow for graceful degradation when the ground truth measurements have position error.

The Earth Mover's Distance, **EMD**, measures the spatial distance between two probability distributions over a region. It was introduced as a spatially robust metric for image matching [162, 187]. Computationally, it is the minimum cost of morphing one distribution into the other. The total cost is the amount of density moved times the distance moved. It can be formulated as a transportation problem [44]. We used the following linear time variant of EMD [162]:

$$\widehat{EMD}(P, Q^D) = \min_{\{f_{ij}\}} \sum_{i,j} f_{ij} d_{ij} + |\sum_i P_i - \sum_j Q_j^D| \max_{i,j} d_{ij}$$

under the constraints:

$$(1) f_{ij} \geq 0 \quad (2) \sum_j f_{ij} \leq P_i \quad (3) \sum_i f_{ij} \leq Q_j^D, \tag{5.6}$$

$$(4) \sum_{i,j} f_{ij} = \min(\sum_i P_i, \sum_j Q_j^D)$$

where each $f_{ij}$ represents the amount of density transported (or the *flow*) from the $i$th supply to the $j$th demand and $d_{ij}$ is the *ground distance* between bin $i$ and bin $j$ in the distribution. Equation 5.6 is therefore attempting to minimize the amount of density movement such that the total density is preserved after the move. Constraint (1) allows transporting density from $P$ to $Q^D$ and not vice versa. Constraint (2) prevents more density to be moved from a location $P_i$ than is there. Constraint (3) prevents more density to be deposited to a location $Q_j^D$ than is there. Constraint (4) is for feasibility: such that the amount of density moved does not exceed the total density found in either $P$ or $Q^D$. Solving this problem requires global optimization on the whole map, making this metric quite computationally expensive.

A larger EMD indicates a larger difference between two distributions while an EMD of zero indicates that two distributions are the same. Generally, saliency maps that spread density over a larger area have larger EMD values (i.e., worse scores) as all the extra density has to be moved to match the ground truth map. EMD penalizes false positives proportionally to the spatial distance they are from the ground truth.

## ■ 5.5 Summary of metric behaviors

This section summarizes the findings of multiple experiments from Bylinskii et al. [36] to discuss some general properties and behaviors of metrics, and provide some guidance for selecting metrics for evaluation.

**Treatment of false positives and negatives:** KL, IG, and SIM penalize models with false negatives significantly more than false positives. If the prediction is close to zero where the ground truth has a non-zero value, the penalties can grow arbitrarily large under these metrics. AUC scores, however, depend on which level sets false positives fall in: false positives in the first level sets are penalized most, but those in the last level set do not have a large impact on performance. Models with many low-valued false positives do not incur large penalties. Therefore, saliency maps that place different amounts of density but at the correct (fixated) locations will receive similar AUC scores. CC and NSS both treat false positives and negatives symmetrically. EMD is least sensitive to uniformly-occurring false negatives because the EMD calculation can redistribute saliency values from nearby pixels to compensate. However, false negatives that are spatially far away from any predicted density are highly penalized. Similarly, EMD's penalty for false positives depends on their spatial distance to ground truth, where false positives close to ground truth locations can be redistributed to those locations at low cost, but distant false positives are highly penalized.

**Relationship between metrics:** Due to their analogous computations, CC and NSS are highly correlated, as are KL and IG. Driven by extreme sensitivity to false negatives, KL, IG, and SIM rank saliency models similarly. Shuffled AUC (sAUC) has low correlations with other metrics because it modifies how predictions at different spatial locations on the image are treated: a model with more central predictions will be ranked lower than a model with more peripheral predictions.

**Selecting metrics for evaluation:** AUC, which is ambivalent to monotonic transformations, has begun to saturate on the MIT Saliency Benchmark and is becoming less capable of discriminating between different saliency models [33]. EMD is computationally expensive to compute and difficult to optimize for. NSS and CC metrics provide the fairest comparison, treating false positives and negatives symmetrically. Closely related mathematically, their rankings of saliency models are highly correlated, and reporting performance using one of them is sufficient. However, under alternative assumptions and definitions of saliency, another choice of metrics may be more appropriate. Specifically, if saliency models are evaluated as probabilistic models, then KL and IG are recommended. Specific tasks and applications may also call for a different choice of metrics. For instance, AUC, KL, and IG are appropriate for detection applications, as they penalize target detection failures. In applications where it is important to evaluate the relative importance of different image regions, such as for image-retargeting, compression, and progressive transmission, metrics like NSS or SIM are a better fit. Overall, we found that either of CC or NSS are a good fall-back option for evaluation, relying on the fewest assumptions about the input. Based on these recommendations, we tend to emphasize evaluation using CC and NSS in this thesis.

In this chapter we discussed the influence of different assumptions on the choice of metrics. We provide code for evaluating and visualizing metric computations (https://github.com/cvzoya/saliency) to add transparency to model evaluation and allow researchers to visualize aspects of saliency models driving or hurting performance.

# Part II

# State of Computational Models of Attention

# Chapter 6

# Where should saliency models look next?

*This chapter is based on: Bylinskii, Z., Recasens, A., Borji, A., Oliva, A., Torralba, A., Durand, F. "Where should saliency models look next?" [European Conference on Computer Vision 2016]*

**W**HERE human observers look in images can provide important clues to human image understanding: where the main focus of the image is, where an action or event is happening in an image, and who the main participants are. The collection of human eye movements can help highlight image regions of interest to human observers, and models can be designed to make computational predictions. The field of saliency estimation has moved beyond the modeling of low-level visual attention to the prediction of human eye fixations on images. This transition has been driven in part by large datasets and benchmarks of human eye movements.



**Figure 6.1.** Recent progress in saliency modeling has significantly driven up performance scores on saliency benchmarks. On first glance, model detections of regions of interest in an image appear to approach ground truth human eye fixations (Fig. 6.3). A finer-grained analysis can reveal where models can still make significant improvements. High-density regions of human fixations are marked in yellow, and show that models continue to miss these semantically-meaningful elements.

For a long while, the prediction scores of saliency models have increased at a stable rate. The recent couple of years have seen tremendous improvements on well-established saliency benchmark datasets [31]. These improvements can be attributed to the resurgence of neural networks in the computer vision community, and the application of deep architectures to saliency estimation. As a result, a large number of neural network based saliency models have emerged in a short period of time, creating a large gap in performance relative to traditional saliency models that are based on hand-crafted features, and learning-based models that integrate low-level features with object detectors and scene context [19, 106, 111, 223, 230]. Neural network-based models are trained to predict saliency in a single end-to-end manner, combining feature extraction, feature integration, and saliency value prediction (Chapter 3.3).

These recent advances in the state-of-the-art and the corresponding saturation of some evaluation scores motivate the questions: Have saliency models begun to converge on human performance and is saliency a solved problem? In this chapter we provide explanations of what saliency models are still missing, in order to match the key image regions attended to by human observers. We argue that to continue to approach human-level performance, saliency models will need to discover increasingly higher-level concepts in images: text, objects of gaze and action, locations of motion, and expected locations of people in images. Moreover, they will need to reason about the relative importance of image regions, such as focusing on the most important person in the room or the most informative sign on the road. In other words, more accurately predicting where people look in images will require higher-level image understanding: **moving beyond image saliency to image importance**. In this chapter, we examine the kinds of problems that remain and what will be required to push performance forward.

## ■ 6.1 Evaluating progress of saliency models

We perform our evaluation on two datasets from the well-established MIT Saliency Benchmark [31]. We use the data from this benchmark because it has the most comprehensive set of traditional and deep saliency models evaluated. The **MIT300** dataset [107] is composed of 300 images from Flickr Creative Commons and personal collections. It is a difficult dataset for saliency models, as images are highly varied and natural. Fixations of 39 observers have been collected on this dataset, leading to fairly robust ground-truth to test models against. The **CAT2000** dataset [20] is composed of 2000 images from 20 different categories, varying from natural indoor and outdoor scenes to artificial stimuli like patterns and sketches. Images in this dataset come from search engines and computer vision datasets [222, 226]. The test portion of this dataset, used for evaluation, contains the fixations of 24 observers.

As of March 2016, of the top 10 (out of 57) models evaluated on MIT300, neural network models filled 6 spots (and the top 3 ranks) according to many metrics[1]. DeepFix [125] and SALICON [103], both neural network models, hold the top 2 spots. The

---

[1]As of April 2018, 15 of the top 15 (out of 77) models on MIT300 are neural networks.

| Saliency model | AUC ↑ | sAUC ↑ | NSS ↑ | CC ↑ | KL ↓ | EMD ↓ | SIM↑ | IG ↑ |
|---|---|---|---|---|---|---|---|---|
| Human limit | 0.92 | 0.81 | 3.29 | 1 | 0 | 0 | 1 | 1.80 |
| DeepFix [125] | **0.87** | 0.71 | **2.26** | **0.78** | 0.63 | **2.04** | **0.67** | 0.67 |
| SALICON [103] | **0.87** | **0.74** | 2.12 | 0.74 | **0.54** | 2.62 | 0.60 | **0.71** |
| BMS [227] | 0.83 | 0.65 | 1.41 | 0.55 | 0.81 | 3.35 | 0.51 | 0.22 |
| IttiKoch[2] | 0.75 | 0.63 | 0.97 | 0.37 | 1.03 | 4.26 | 0.44 | -0.15 |
| Chance | 0.50 | 0.50 | 0 | 0 | 2.09 | 6.35 | 0.33 | -1.67 |

**Table 6.1.** Scores of top-performing neural network models (DeepFix, SALICON) and best non-neural network model (BMS) on MIT300 Benchmark. Top scores are bolded. Lower scores for KL and EMD are better. There has been significant progress since the traditional bottom-up IttiKoch model, but a gap remains to reach human-level performance. Chance and human limit values have been taken from [31, 36].

CAT2000 dataset, a recent addition to the MIT benchmark, has 19 models evaluated to date. DeepFix is the best model on the CAT2000 dataset overall and on all 20 image categories. BMS (Boolean map based saliency) [227] is the best-performing non neural network model across both datasets.

A finer-grained analysis on MIT300 showed that on a per-image level, DeepFix and SALICON alternate in providing the best prediction for ground-truth fixations. In the rest of the chapter, our analyses are carried out on these models. Performances of these models on the MIT benchmark according to the benchmark metrics are provided in Table 6.1 (see Chapter 5 for a description of the metrics).

To begin to explore where these recent large gains in performance are coming from, we visualize the most representative dataset images in Fig. 6.2. These representative images were chosen using a correlation-based greedy approach [78]. We greedily select one image at a time from the MIT300 dataset to best approximate the model score ranking on the MIT Saliency Benchmark, while keeping model performances on the images selected as uncorrelated as possible (to increase diversity of the subset of images selected). We select a subset of $k$ images by optimizing:

$$\text{CSF} = \frac{k\overline{r_{st}}}{\sqrt{k + (k-1)\overline{r_{ss}}}}$$

$$\text{where } \overline{r_{st}} = \frac{1}{k}\sum_{i=1}^{k} corr(s_i, t)$$

$$\overline{r_{ss}} = \frac{k(k-1)}{2}\sum_{i=1}^{k}\sum_{j\neq i}^{k} corr(s_i, s_j)$$

where $s_i$ is a vector of the NSS scores for all models on image $i$, and $t$ is a vector of the NSS scores for all models averaged over all 300 images of the MIT benchmark.

---

[2]Implementation from http://www.vision.caltech.edu/~harel/share/gbvs.php.

**Figure 6.2.** Saliency model ranking is preserved when evaluating models on this subset of 10 images as when evaluating them on the whole 300-image benchmark. These images help to accentuate differences in model performance. These images contain people at varying scales, as well as text (small here) amidst distracting textures.

We find that a subset of $k = 10$ images can already rank the saliency models on the MIT benchmark with a Spearman correlation of 0.97 relative to their ranking on all dataset images. These images help to accentuate differences in model performance. By visualizing the predictions of some of the top and bottom models on these images (Fig. 6.3), we can see that driving performance is a model's ability to detect people and text in images in the presence of clutter, texture, and potentially misleading low-level pop-out.



**Figure 6.3.** Some of the best and worst model predictions on a few of the representative images from Fig. 6.2. Unlike traditional bottom-up models, recent neural network models can discover faces, text, and object-like features in images, prioritizing them over textures and low-level features appropriately, to better approximate human fixations.

**Figure 6.4.** Two types of Mechanical Turk tasks were used for gathering annotations for the highly-fixated regions in an image. These annotations were then used to quantify where people look in images.

## ■ 6.2  Quantifying where people and models look in images

To understand where models might fail, we must first understand where people look. Our goal is to name all the image regions lying beneath the high-density locations in fixation heatmaps. We computed fixation heatmaps aggregated over all observers on an image (39 observers in the MIT300 dataset, for a robust ground truth). Then we thresholded these ground truth heatmaps at the 95th percentile and collected all the connected components. This produced an average of 1-3 regions per image for a total of 651 regions.

The resulting region outlines were plotted on top of the original images and shown to Amazon Mechanical Turk (MTurk) participants with the task of selecting the labels that most clearly describe the image content that people look at (Fig. 6.4a). The labels provided for this task were not meant to serve as an exhaustive list of all objects, but to have good coverage of label types, with sufficient instances per label. If an image contained multiple image regions, only one would be displayed to participants at a time. Participants could select out of 15 different label categories as many labels as were appropriate to describe a region. For each image region, we collected labels from a total of 20 participants. Majority vote was used to assign labels to regions. A region could have multiple labels in case of ties. For further analyses, related labels (e.g. "animal face", "part of an animal", etc.) were aggregated to have sufficient instances per label type. Not all regions are easily nameable, and in these cases participants could select the "background" or "other" labels. To account for these image regions to which simple labels could not be assigned, a second question-based MTurk task was deployed, described in the next section. Table 6.2 summarizes the labels assigned to the image regions frequently fixated by human viewers.

**Table 6.2.** What do people look at in images? Regions in images receiving a high density of eye fixations were labeled by MTurk participants. We summarize the 681 labels assigned to the 651 regions by MTurk participants.

| Region type | Number of instances |
|---|---|
| Object | 264 |
| Part of a person | 97 |
| Legible text | 84 |
| Part of a face | 67 |
| Part of an animal | 42 |
| Crowd of people | 33 |
| Face | 27 |
| Other | 19 |
| Person | 14 |
| Background | 13 |
| Animal face | 6 |
| Illegible text | 5 |
| Head | 5 |
| Non-english text | 3 |
| Symbol | 2 |

### ■ 6.2.1 What do models miss?

Given labels for all the highly-fixated image regions in the MIT300 dataset, we intersected these labeled regions with the saliency maps of different computational models. To determine if saliency models made correct predictions in these regions, we calculated whether the mean saliency in these regions was within the 95-th percentile of the saliency map for the whole image. We then tallied up the types of regions that were most commonly under-predicted by models. In Table 6.3 we provide the error percentages, by region type, where saliency models assigned a value less than the percentile threshold to the corresponding regions. Our analyses are performed over DeepFix and SALICON models on the MIT300 dataset, and on DeepFix on the CAT2000 dataset (additional analyses are provided in the Supplemental Material of [33]). The four categories chosen from the CAT2000 dataset are ones that contain natural images with a variety of objects and settings.

About half the failure modes are due to misdetections of parts of people, faces, animals, and text. Such failure cases can be ameliorated by training models on more instances of faces (partial, blurry, small, non-frontal views, occluded), more instances of text (different sizes and types), and animals. However, the labels "background", "object", and "other" assigned to image regions by MTurk participants originally accounted for about half of model errors on MIT300.

A second MTurk task was designed to better understand the content found in these harder-to-name image regions. Participants were asked to answer binary questions, such as whether or not a highlighted region in an image is an object of gaze or action in the image (see Fig. 6.4b and Table 6.2). The results of this task allowed us to further break down model failure modes, and account for 60% of total mispredictions on MIT300 and

| Dataset | MIT300 | | CAT2000 | | | |
|---|---|---|---|---|---|---|
| Model | DeepFix | SALICON | DeepFix | | | |
| Image category | All | | Social | Action | Indoor | Outdoor |
| Part of main subject | 31% | 36% | 49% | 68% | 12% | 24% |
| Unusual element | 18% | 16% | 33% | 63% | 8% | 8% |
| Location of action/motion | 16% | 16% | 67% | 78% | 8% | 11% |
| Text | 16% | 13% | 6% | 5% | 8% | 29% |
| Part of a person | 15% | 14% | 23% | 37% | 8% | 5% |
| Possible location for a person | 15% | 7% | 6% | 24% | 10% | 11% |
| Object of action | 14% | 15% | 27% | 51% | 0% | 3% |
| Object of gaze | 11% | 11% | 50% | 44% | 0% | 0% |
| Part of a face | 6% | 8% | 46% | 7% | 0% | 0% |
| Part of an animal | 5% | 5% | 3% | 10% | 0% | 0% |
| Other | 40% | 40% | 3% | 2% | 61% | 37% |

**Table 6.3.** Labels for under-predicted regions on MIT300 and CAT2000 datasets. Percentages are computed over 681 labels assigned to 651 regions (some regions have multiple labels so percentages do not add up to 100%). See Fig. 6.6 for visual examples.

39%-98% of mispredictions on four categories of CAT2000. The remaining failure modes (labeled "other") vary from image to image, caused by low-level features, background elements, and other objects or parts of objects that are not the main subjects of the photograph, nor are objects of gaze or action. Later in this chapter, the most common failure modes are explored in greater detail. Examples are provided in Fig. 6.6.

## ■ 6.2.2 What can models gain?

With the region annotations obtained from our MTurk tasks, we performed an analysis complementary to that in Sec. 6.2.1. Instead of computing model misses across different image regions, here we estimate the potential gains to models if specific image regions were correctly predicted. A region is treated as a binary mask for the image, and a modified saliency map is computed as a combination of the original saliency map and ground truth fixation map. For each region type (e.g. "part of a person", "object of gaze"), we compute modified saliency maps. We replace model predictions in those regions with ground truth values obtained from the human fixation map (e.g., Fig. 6.5). Fig. 6.5 provides the score improvements of the modified models on the MIT300 benchmark. This analysis is meant to provide a general sense of the possible improvement if different prediction errors are ameliorated. We include improvements in Normalized Scanpath Saliency (NSS) and Information Gain (IG) scores, which follow the distribution of region types in Table 6.3. We found that the Area under ROC Curve (AUC) metrics have either saturated or are close to saturation. The focus of saliency evaluation should turn instead towards metrics that can continue to differentiate between models, and that can measure model performances at a finer-grained level (Sec. 6.3).

**Figure 6.5.** Improvements of DeepFix and SALICON models on MIT300 if specific regions were accurately predicted. Performance numbers are over all 300 benchmark images, where regions from the ground truth fixation map are substituted into each model's saliency maps to examine the change in performance (top row). The percentage score improvement is computed as a fraction of the score difference between the original model score and the human limit (from Table 6.1).



**Figure 6.6.** Regions often fixated by humans but missed by computational models.

**Figure 6.7.** Saliency prediction failure cases for faces: (a) Face saliency is underestimated when faces are small, non-frontal, or not centered in an image; (b) Sometimes the actions in a scene are more salient to human observers than the participants, but saliency models can overestimate the relative saliency of the faces; (c) Face detection can fail on depictions (such as in posters and photographs within the input images) which often lack the context of a body, or appear at an unusual location in the image.

## ■ 6.2.3  The importance of people

A significant part of the regions missed by saliency models involve people (Table 6.3): people within the salient region, or people acting on or looking at a salient object. In this section we provide a deeper analysis of the images containing people. To expand our analysis, we annotated all the people's faces in the MIT300 images with bounding boxes. This provided a more complete set of annotations than the regions extracted for the MTurk labeling tasks, where only the top 1-3 most highly-fixated regions per image were labeled. In this section we compute the importance of faces in an image following the approach of Jiang et al. [103]: given a bounding box for an object in an image, the maximum saliency value falling within the object's outline is taken as the object's **importance score** (the maximum is a good choice for such analyses as it does not scale with object size). This will be used to analyze if saliency models are able to capture the relative importance of people in scenes.

Across the images in MIT300 containing only one face (53 images), the face is the most highly fixated region in 66% of the images, and the DeepFix model correctly predicts this in 77% of these cases. Out of the 53 images with faces, the saliency of the face is underestimated (by more than 10% of the range of saliency values) by the DeepFix model in 15 cases, and overestimated in 3 cases. In other words, across these images, the DeepFix model does not assign the correct relative importance to the face relative to the rest of the image elements in a third of the total cases. Some of these examples are provided in Fig. 6.7. Note that the importance of faces extends to depictions of faces as well: portraits or posters containing human faces in images. Human attention is drawn to these regions, but models tend to miss these faces, perhaps because they are lacking the necessary context to discover them.

**Figure 6.8.** Although recent saliency models have begun to detect faces in images with high precision, they do not assign the correct relative importance to different faces in an image. This requires an understanding of the interactions in an image: who is participating in an action and who has authority. Facial expressions, accessories/embellishments, facial orientation, and position in a photo also contribute to the importance of individual faces. We assign an importance score to each face in an image using the maximum ground truth (fixation) or predicted (saliency) density in the face bounding box. These importance scores, ranging from 0 to 1, are included above each bounding box.

Similarly to the analysis in Sec. 6.2.2, here we quantify the performance boost of saliency models if the saliency of faces were always correctly predicted. We used the same procedure: to create the modified saliency map for an image we assign the ground truth saliency value to the bounding box region and the predicted output of the model to the remaining part of the image. The DeepFix model's Normalized Scanpath Saliency (NSS) score on the MIT300 benchmark improves by 7.8% of the total remaining gap between the original model scores and human limit, when adding ground truth in the face bounding boxes. Information Gain (IG) also goes up 1.8%. Improving the ability of models to detect and assign correct relative importance to faces in images can provide better predictions of human eye fixations.

## ■ 6.2.4  Not all people in an image are equally important

Considering images containing multiple faces, we measure the extent to which the computational prediction of the relative importance of the different faces matches human ground-truth fixations. For all the faces labeled in an image, we use the human fixation maps to compute the importance score for each face, and analogously we use the saliency map to assign a predicted importance score to the same faces. Since both fixation and saliency maps are normalized, each face in an image will receive an importance score ranging from 0 to 1. A score of 1 occurs when the face bounding box overlapped a region of maximum density in the corresponding fixation/saliency map. Interpreted

**Figure 6.9.**  Example images containing text that receive many fixations by human observers, but whose saliency is under-estimated by computational models.  Text labels can be used to give the observer more information. For instance, the description of a warning or a book are more informative to observers than the warning or book title itself.  These regions receive more eye fixations.  The informativeness of text also depends on the context of the observer: most observers fixated the only piece of English text on the box of chocolates.

in terms of ground truth, this is the face that received the most fixations.

Across the images with more than one visible face, the average Spearman correlation between the ground truth and predicted face importance values is 0.53.  This means that for many images, the relative ordering assigned by the saliency model to people does not match the importance given by human fixations.  As depicted in Fig. 6.8, discovering the most important person in the image is a task that requires higher-level image understanding.  Human participants tend to fixate people in an image that are central to a depicted action, a conversation, or an event; people who stand out from the crowd (based on some high-level features like facial expression, age, accessories, etc.).

■ **6.2.5  The informativeness of text**

In the MIT300 and CAT2000 datasets, most text, large or small, has attracted many human fixations, with regions containing text accounting for 7% of all highly-fixated image regions.  While text has been previously noted as attracting human visual attention [38], not all text is equal.  The informativeness of text in the context of the

rest of the image, or the interestingness of the text on its own can affect how long individual observers fixate it, and what proportion of observers look at it. There are thus a number of reasons why the human ground truth might have a high saliency on a particular piece of text, and some of those reasons depend on understanding the text itself - something that computational models currently lack (Fig. 6.9).

To expand our analysis on text regions, we annotated all instances of text present in the MIT300 dataset with bounding boxes. The DeepFix model's NSS scores improves by 7.8% of the total remaining gap between the original model scores and human upper bound, when adding ground truth in the text bounding boxes. Its IG score improves by 4.4%. Overall, an accurate understanding of text is another step towards better predictions of human fixations.

## ■ 6.2.6  Objects of gaze and action

Another common source of missed predictions are objects of gaze and/or action. These are objects or, more generally, regions in an image that are looked at or interacted with by one or more persons in an image. In Fig. 6.10, we include 4 images from the MIT300 dataset that include objects of gaze missed by both DeepFix and SALICON. Training saliency models to explicitly follow gaze can improve their predictive power of modeling the saliency of the entire scene [201]. In the last column of Fig. 6.10 we show the predictions that can be made possible by a computational model specifically trained to predict gaze [176]. For each person in an image, this model predicts the scene saliency from the vantage point of the individual selected. Its output consist of a heat map representing a combination among the different gaze predictions; that is, a map highlighting the objects people are looking at in the image. The procedure we used to build the final gaze maps is described below.

1. Using face bounding boxes we compute the output of the model for each of the people in the image.

2. Using the importance score for each of the people in the image, we can weight each of the gaze maps by its relative importance in the image.

3. Adding up all the weighted maps we compute the final output map. These output maps provide a distribution over where each of the people is looking.

The final weighted map captures the objects where people are looking and their relative importance to the full image. The gaze-following model only works when gaze information can be extracted from the orientation of the head and, if visible, the location and orientation of the eyes. However, the orientation of the body and location of body parts (specifically the hands) can provide additional clues as to which objects in an image are relevant from the vantage point of different people in the image, even if not fully visible. Detecting such objects of action remains a problem area for saliency models (Fig. 6.11).

**Figure 6.10.** Both top neural network saliency models perform worse on these images than on any other images in the MIT300 dataset labeled with objects of gaze. The yellow outlines highlight high-density regions in the ground truth fixation map that were labeled by MTurk participants as regions on which the gaze of someone in the image falls. A model that explicitly predicts the gaze of individuals in an image can locate these objects of gaze [176]. The last row is a failure of the gaze-following model, requiring an understanding of actions that is beyond just gaze.

**Figure 6.11.** Images in the MIT300 dataset labeled to contain objects of action - i.e. objects being acted on, or interacted with, by people in the scene. Included are images where both deep learning saliency models, DeepFix and SALICON, underestimate the saliency of these regions. Notice that the significance of these objects can not be inferred from gaze information (unlike in Fig. 6.10), since in all of these cases no one in the image is looking at the objects of interest. The yellow outlines highlight high-density regions in the ground truth fixation map that were labeled by MTurk participants as regions containing objects being acted on.

**Figure 6.12.** A finer-grained test for saliency models: determining the relative importance of different sub-images in a panel. (a) A panel image from the MIT300 dataset. (b) The saliency map predictions given the panel as an input image. (c) The maximum response of each saliency model on each subimage is visualized (as an importance matrix).

## ◾ 6.3  Future directions

As the number of saliency models grows and score differences between models shrink, evaluation procedures should be adjusted to elucidate differences between models and human eye movements. This calls for finer-grained evaluation metrics, datasets, and prediction tasks. Models continue to under-predict crucial image regions containing people, actions, and text. These are precisely the regions with greatest semantic importance in an image, and become essential for saliency applications like image compression and image captioning. Aggregating model scores over all image regions and large image collections conceals these errors. Moreover, traditionally favored saliency evaluation metrics like the AUC can not distinguish between cases where models predict different relative importance values for different regions of an image. As models continue to improve in detection performance, measuring the relative values they assign to the detected objects is the next step. This can be accomplished with metrics like the Normalized Scanpath Saliency (NSS) and Information Gain (IG), which take into account the range of saliency map values during evaluation [36]. Finer-grained tasks like comparing the relative importance of image regions in a collection or in a panel such as the one in Fig. 6.12 can further differentiate model performances. Finer-grained datasets like CAT2000 [20] can help measure model performance per image type.

Recent saliency models with deep architectures have shown immense progress on saliency benchmarks, with a wide performance gap relative to previous state-of-the-art. In this chapter we demonstrated that a finer-grained analysis of the top-performing models on the MIT Saliency Benchmark can uncover areas for further improvement to narrow the remaining gap to human ground truth[3].

---

[3]Additional results are provided in the Supplemental Material of [33].

# Part III

# Crowdsourcing Human Attention

# Chapter 7

# BubbleView: an interface for crowdsourcing image importance

*This chapter is based on: Kim, N.W.\*, Bylinskii, Z.\*, Borkin, M.A., Gajos, K., Oliva, A., Durand, F., Pfister, H. "BubbleView: an interface for crowdsourcing image importance maps and tracking visual attention" [Transactions on Computer-Human Interaction 2017]*

**E**YE tracking has proven useful for studying the cognitive processes involved in visual information processing, including which visual elements people look at first and spend the most time on [99, 143] (Chapter 3). Eye tracking is widely used for conducting usability studies for human-computer interfaces [99, 151], for designing gaze-based and attention-aware user interfaces [142, 143] or for collecting gaze data to build saliency prediction models [20, 107].

Commercial eye-trackers mostly use specialized hardware such as advanced infrared sensors and high-quality cameras to accurately track eye positions and movements [1]. However, they often require high-cost equipment and invasive calibrations (e.g., Eye-Link, ISCAN), which means it is difficult to scale to large scale studies beyond controlled lab environments. Recent appearance-based methods attempt to address this issue by enabling eye tracking on affordable cameras built into personal devices [89, 123, 224] (Chapter 3.2). However, these methods have not yet seen widespread adoption, as they still suffer in accuracy and robustness, and impose set-up constraints (camera quality, lighting conditions).

On the other hand, cursor-based attention tracking is based on the correlation between gazes and cursor locations [74, 88, 184] and reduce the need to handle variations in real-world settings in camera-based methods; e.g., calibrations, ambient lighting, etc. The most popular cursor-based approach uses a moving window continuously following the position of the cursor to reveal a portion of the screen in normal resolution [101] (Chapter 3.1).

Our **BubbleView** methodology is a cursor-based, moving-window approach to collect clicks on static images as a proxy for eye fixations. BubbleView presents blurred images and allows participants to click around to reveal small circular "bubble" regions

(a) saliency map



(b) importance map



(c) eye-tracking set-up



(d) BubbleView set-up

**Figure 7.1.**    Just as the pattern of human eye fixations can be used as a heatmap of saliency for an image (a), the pattern of BubbleView clicks can be used as a heatmap of importance for an image (b). An eye tracking set-up (pictured: EyeLink1000) is a way to collect human eye fixations in the lab setting (c), whereas the BubbleView interface can be launched online and feasibly scale up the collection of crowdsourced data (d).

of the image at the original resolution (Figure 7.1). This is intended to loosely approximate a blurred periphery and the confined area of focus of the human eye fovea (Chapter 5.1).

Compared to natural viewing, BubbleView and related cursor-based methodologies slow down the exploration patterns of participants, because choosing where to move the mouse and click is a slower cognitive process than moving eyes around an image. Because of this, we refer to the pattern of BubbleView clicks on an image as the **importance map** for the image. We intend for importance to encapsulate image regions that are not only more attention grabbing initially (salient), but also regions that people spend more time on because they are more relevant, or interesting, to the task-at-hand.

BubbleView is especially well suited to capturing image regions of most importance when a directed task is provided (as compared to free viewing). Our initial target setting, first presented in Kim et al. [112] was to show that BubbleView clicks can provide

a good approximation for eye movements when participants are asked to describe the content of information visualizations (graphs, charts, tables). In Borkin et al. [24], we further showed that knowing where people look can provide clues about what they store in memory and recall about an information visualization. Like eye tracking, BubbleView can provide important insights about human perception and cognition, but at a lower data collection cost than eye tracking. It can easily scale up data collection to many participants and images, and be launched remotely to enable online crowdsourcing.

In this chapter, we validate that BubbleView generalizes to approximating eye fixations on different image types and under different task constraints. Specifically, we show that:

- BubbleView clicks can successfully approximate eye fixations on information visualizations, natural images, and websites, in both a free-viewing condition and with a description task;

- Compared to related methodologies based on a moving-window approach [103], BubbleView clicks provide more reliable and less noisy data;

- The number of BubbleView clicks in different image regions can be used to measure the relative importance of those image regions.

We present the BubbleView methodology with the interested experimenter in mind who may consider it for crowdsourcing an experiment, or for an evaluation that would typically be conducted with an eye tracker in a conventional laboratory setting. While prior work contains some initial validation that a cursor-based interface can serve as a proxy for eye tracking [10, 103, 112], we conducted an extensive quantitative analysis by running 10 experiments with 28 different parameter combinations, on Amazon's Mechanical Turk. Our experiments were carried out on 5 different datasets, spanning information visualizations [24], natural images [103, 223], static webpages [194], and graphic designs [154]. We varied task type (free-viewing, describing) and task duration, image blur kernel, and bubble radius. We compared BubbleView clicks not only to eye fixations [24, 223], but also to mouse movements [103], and to explicit importance annotations [154]. Our contributions include:

1. The BubbleView interface which can be launched online for the cheap, feasible collection of crowdsourced data, provided at `massvis.mit.edu/bubbleview`;

2. A thorough analysis of how different experimental parameters affect BubbleView click data, and guidelines about how to choose an appropriate setting of parameters for a given experiment;

3. A discussion of how BubbleView can be used to approximate eye fixations collected in a controlled lab setting;

4. A proposed list of applications of the BubbleView methodology, including for the measurement of image importance, image-based question-answering tasks, and training computational models of saliency/importance.

## ■ 7.1 Designing experiments with BubbleView

The BubbleView methodology is intended to approximate a blurred periphery, and users click on images to reveal small, circular regions ("bubbles") at the original resolution (Figure 7.2). This is similar to having a confined area of focus like the eye fovea (Chapter 5.1). Different blur levels and bubble sizes can be used to approximate different eye tracking setups, with different visual angles (Figure 7.14).

In comparison to the moving-window approach which records continuous mouse movements, our approach records discrete mouse clicks where each click represents a conscious choice made by the user to reveal a portion of the image. As the clicks correspond to individual points of interest, we directly compare them to eye fixations.



**Figure 7.2.** Two different versions of the BubbleView interface for two task types, for gathering task-based (a) and task-free (b) clicks, as approximations to similar eye tracking experiments.

*Tasks and image types for attention experiments*
We evaluated our BubbleView interface on two tasks: free-viewing and description, and four image types: natural scenes, information visualizations, static webpages, and graphic designs. Here we discuss the motivations behind these design choices.

During **free-viewing**, participants are not given a task but are instructed to freely look around the image. Free-viewing is commonly used in eye tracking experiments to study the human perception of natural scenes, because it can avoid large task-dependent effects. It is often assumed the eyes are drawn to conspicuous image elements, and attention proceeds in a bottom-up manner, guided by the image features rather than a high-level task[1]. This assumption has motivated the use of free-viewing for collecting ground truth data for saliency datasets [118], where the pattern of eye fixations can be interpreted as the **saliency map** for the image (Figure 7.1). Most saliency datasets

---

[1]Alternative views posit that free-viewing is not task-free, but permits participants to choose their own internal agendas/tasks [161, 207, 208]. Even under this interpretation, averaging data over many participants, each of which may have their own agenda, has the effect of averaging out the task and providing an approximately task-independent aggregate measurement.

have been collected using free-viewing[2].  Computational models are in turn trained and tested on saliency datasets as a proxy for human attention.  We are similarly motivated by the computational applications that can be built by training models on large attention datasets (e.g., [35]).

Compared to natural viewing, cursor-based moving-window methodologies naturally slow down visual exploration patterns.  By providing a cognitively-demanding task, these exploration patterns can be slowed down further to bring more intentionality to each click.  In the **description** task, participants are required to type a description of the image while using the BubbleView interface to explore the image.  The descriptions naturally depend on the image regions clicked on.  This task is well suited to images with an underlying message or concept that needs careful examination to decipher.  We used the description task with visualization images from the MASSVIS [24] dataset[3], and website images from the FiWI [194] dataset.  We also tested the free-viewing task with the FiWI images, because the eye-tracking data from this dataset was collected with free-viewing, and we wanted to approximate the original experiment.  For the same reason, we ran the free-viewing task on natural images from the OSIE [223] dataset.  For the graphic designs in the GDI dataset [154], which have importance annotations rather than eye fixations, we chose a free-viewing task.  We chose this task because we found that the graphic designs could not be easily summarized by a description (i.e., some images required further context, not all were English, some had few visual elements, etc.).

Tasks deviating from description and free-viewing are beyond the scope of this chapter, although they are common in user interface research [11, 43, 67, 99].  For instance, for testing websites or application interfaces, participants may be asked to perform tasks such as searching for a particular element or option, navigating to a particular region of the image or page, or answering questions.  Related moving-window methodologies have previously been validated in the context of web navigation, program debugging, and question-answering [10, 101, 132, 192, 205].  These tasks can be quite specific to the interface being evaluated.  We used two task types that can generalize (without modification) to a large collection of different image types.  Our BubbleView tool is available to the research community so future work can investigate the generalizability of this tool for other tasks.

*Implementation*
We implemented a web-based BubbleView interface that takes a directory of images as input and displays a subset of the images in random sequence, blurring each one. Participants receive a set of task instructions and can click to reveal bubble regions (Figure 7.2).  A demo is available at massvis.mit.edu/bubbleview.

---

[2]A list of eye tracking datasets and their attributes is available at: http://saliency.mit.edu/datasets.html

[3]In the MASSVIS eye-tracking set-up participants also provided image descriptions, but they did so at the end, not during, the viewing session.  This is because memorability was part of the original study, whereas it is not here.

The experimenter has a choice of parameters:

- **Task type:** the instructions given to participants. We used two different versions of the interface for a description task with an input text field (Figure 7.2a), and a free-viewing task with no additional inputs from participants (Figure 7.2b). Alternative tasks are possible.

- **Time:** the viewing time per image, which depends on the task. For the description task, we did not constrain the time. For the free-viewing task, we fixed time per image to be either 10 or 30 seconds, depending on the experiment.

- **Blur sigma:** the size of the Gaussian blur kernel (in pixels) to apply to each image to mimic peripheral vision. This is a fixed quantity over the whole image, and is constant across all images in the sequence. In our studies, we manually selected a blur value per image dataset to distort image text beyond recognition. We wanted the level of detail to be sufficient for reading only within regions of focus.

- **Bubble radius:** the size of the focus area (in pixels) that is deblurred during a click to mimic foveal vision. In our studies, we varied this size depending on other task constraints, but often stayed within 1-2 degrees of visual angle of the eye tracking setups used for the ground-truth eye movement datasets.

- **Mouse modality:** although we originally designed BubbleView for collecting mouse clicks, we extended it to allow bubble regions to be exposed during continuous mouse movements (as in Jiang et al. [103]). We discuss the differences between the two modalities in Section 7.6.

The experimenter may also choose the number of images displayed in a sequence. In our description task, participants were able to continue to the next image after writing a minimum number of characters (150 in our experiments). In the free-viewing task, once the fixed time per image elapsed, the next image in the sequence was presented.

We also developed a monitoring interface to inspect experimental results (Figure 7.3). The purpose of the interface is to take a quick glance at the bubbles collected, before the main analysis. For each image, the experimenter can see the bubbles and (if applicable) text descriptions generated by each participant. Adjusting the slider allows exploration of the temporal sequence and evolution of bubble clicks and description text over time. The experimenter can also see how the blurred image looked to the participant to investigate why a region may have been clicked. This interface can be used to check if an experiment is running as intended in real time.

## ■ 7.2 Analysis overview

Across all the experiments comparing BubbleView clicks to eye fixations we used the same set of analyses, which we describe here. We compared how well the distribution

**Figure 7.3.** Monitoring interface for manually inspecting the results of experiments. An experimenter can use a slider to explore the temporal sequence and evolution of bubble clicks and description text, for each image and participant.

of BubbleView clicks approximates the distribution of eye fixations. Given a set of eye fixations on an image, we generate a fixation map by blurring the fixation locations with a Gaussian, with a sigma equal to one degree of the visual angle. Similarly, given a set of BubbleView mouse clicks on an image, we compute a **BubbleView click map** by blurring the click locations with a Gaussian with the same sigma as for the ground truth fixation maps. We used a sigma of 10 for the OSIE dataset, and a sigma of 25 for the MASSVIS and FiWI datasets. More generally, we refer to both fixation and click maps in this paper as **importance maps** for an image.

For evaluation, we use the two metrics found to give the most un-biased evaluation in Chapter 5: Pearson's Correlation Coefficient (CC) and Normalized Scanpath Saliency (NSS). While the two metrics provide complementary evidence for our conclusions, the NSS metric also allows us to account for differences in attentional consistency between participants (inter-observer congruency) across datasets. Specifically, if different eye tracking participants look at different regions of the image, they can not be used to predict each other's fixations. In these cases, BubbleView clicks will also not be as predictive of the fixations. For a fair evaluation, we normalize the BubbleView scores by the consistency of the eye tracking participants in a given dataset.

Consistency between eye tracking participants is measured in the following way: the fixations of all but one observer (i.e., N-1 observers) are aggregated into a fixation

map which is used to predict the fixations of the remaining observer. This is repeated by leaving out one observer at a time, and then averaging the prediction performance to obtain the resulting inter-observer congruency (**IOC**) or inter-subject consistency [18, 147, 219]. We measure IOC using the NSS metric.

We first compute the NSS score of the BubbleView click map at predicting all the eye fixations collected on an image, across all the observers. Then we normalize this score by the IOC of the eye tracking participants on that dataset. The resulting **normalized NSS** score can be interpreted as: the percent of the eye fixations accounted for, or predicted by, the BubbleView clicks.

We also consider performance when the number of study participants is taken to the limit, to get an upper bound on performance and determine if any systematic differences exist between methodologies that can not be reduced by gathering more data. To do this, we measure the ability of BubbleView click maps to predict ground-truth fixation locations, for different numbers of BubbleView participants. We obtain an NSS score for different numbers of participants $n$, by randomly selecting $n$ participants for each of 10 splits, and averaging the results. Then we fit these scores to the power function $f(n) = a * n^b + c$, constraining $b$ to be negative. Taking $n$ to the limit, $c$ is the NSS score at the limit. In cases where the total number of BubbleView participants for a particular experiment is not enough for a robust model fitting, we omit this analysis.

## ■ 7.3 Experiments comparing BubbleView clicks to eye fixations

## ■ 7.3.1 Experiment 1: comparison to eye fixations on information visualizations

We began by exploring how well BubbleView clicks on information visualizations gathered on MTurk approximate eye fixations collected in a controlled lab setting. In the initial experiments in Kim et al. [112], we had gathered BubbleView data on 51 visualizations with a bubble radius size of 16 pixels. Here we extended these experiments to explore the effect of bubble radius size and number of participants on the quality of BubbleView data. We varied the bubble radius between 16 and 40 pixels, and collected up to 40 participants worth of clicks per image.

> *Motivating questions*
> - How does bubble radius size affect performance?
>
> - How many BubbleView participants is enough?

*Stimuli*
The MASSVIS dataset contains over 5,000 information visualization images, of which 393 "target" images contain the eye movements of 33 participants free-viewing each image for 10 seconds as part of a memory test at the end of the study [24]. In the eye tracking set-up, images were shown full-screen with a maximum dimension of 1000 pixels

**Figure 7.4.** Example images from the MASSVIS dataset. Dataset images (a), with corresponding ground-truth fixation maps (b) and BubbleView click maps (c). We show cases where BubbleView maps have high similarity, and cases with low similarity, to fixation maps.

to a side, where 1 degree of viewing angle corresponded to 32.6 pixels. Participants made on average 39 fixations per image, or 3.9 fixations/sec.

We selected 202 from the total 393 target images, spanning infographic, news media, and government publication categories (Figure 7.4). We chose visualizations that had sufficiently large text and enough context to understand them without requiring specialized knowledge. We resized the images to half their original size with a maximum dimension of 500 pixels to a side. The images were blurred with a sigma of 40 pixels, which we found distorted the text in these images beyond legibility [24, 112].

*Method*
We ran a series of experiments to progressively find a bubble radius that best approximates eye fixations: **Exp. 1.1** with one set of 51 images and bubble radius sizes of 16, 24, and 32 pixels respectively, **Exp. 1.2** with another set of 51 images and bubble radius sizes of 24, 32, and 40 pixels, and **Exp. 1.3** with the remaining 100 images with a bubble radius of 32 pixels, which we determined from the first two experiments to produce good data quality. A bubble radius of 32 pixels corresponds to about 2 degrees of visual angle in the eye tracking studies on the original-sized images.

In a single HIT, participants were shown a random sequence of 3 images, and asked to describe each image with no time constraints on the task, allowing for individual differences in the time to write image descriptions.

For Exp. 1.1, we requested enough HITs so that each image would be seen by an

average of 40 participants. From this experiment we found that 10–15 participants are sufficient for achieving high similarity scores to eye fixations, and proceeded to collect an average of 10–15 participants for each image in Exp. 1.2 and Exp. 1.3.

*Results on bubble size*

Participants explored each image for an average of 3 minutes, iterating between clicking around and typing text. As bubble size increased, the number of clicks and total task time monotonically decreased. Participants made an average of 103 clicks per image (0.5 clicks/sec) with a bubble radius of 16 pixels, 65 clicks (0.3–0.4 clicks/sec) with a bubble radius of 32, and 55 clicks (0.3 clicks/sec) with a bubble radius of 40 pixels. Depending on the bubble size, participants spent 15–30% of the task time clicking, and the rest of the time typing a description. After receiving a number of participant complaints about task difficulty at a bubble radius of 16 pixels, we discontinued the use of this bubble radius in future experiments.

We computed the similarity between the BubbleView click maps and ground truth fixation maps across all images for all settings of bubble radius (Table 7.3.1). To make scores comparable, we set the number of participants $n = 10$ when computing the BubbleView click maps (the common denominator across all experiments). The similarities between the BubbleView click maps and the fixation maps were close across all bubble radius sizes (CC = 0.82–0.86). Because the different subsets of the MASSVIS dataset used in Exp. 1.1–1.3 had different inter-observer consistency (IOC) values[4], normalized NSS scores are more comparable across experiments than raw NSS scores. The normalized NSS score was very similar across all bubble radius sizes, with BubbleView clicks accounting for an average of 89–90% of eye fixations with 10 participants, and climbing up to 92% for larger numbers of participants ($n \geq 18$). Running a one-way ANOVA with bubble size as the factor, we did not find any significant effects of radius size on the similarity of clicks to fixations, under either of CC and NSS scores ($F < 1$ for all comparisons). Although the number of clicks changed, the overall pattern of bubble clicks remained the same (Figure 7.5).

> *Take-aways:* no significant differences were found between bubble sizes in terms of similarity of BubbleView clicks to eye fixations. Bubble sizes in the range 24–40 pixels were found appropriate. Smaller bubble sizes increased the task time and effort.

*Results on number of participants*

In Exp. 1.1 we collected an average of 40 participants of BubbleView clicks per image to investigate how BubbleView maps change with the number of participants (Figure 7.6). As described in Section 7.2, we fit power functions to the NSS scores for different numbers of participants to extrapolate performance. We found that after about 10–15

---

[4]This is an artifact of the images being different in the different subsets. In particular, Exp. 1.1 ended up containing more news media images and less government and infographic images than Exp. 1.2.

(a) input image     (b) bubble radius = 16     (c) bubble radius = 24     (d) bubble radius = 32

**Figure 7.5.** We found few differences in the resulting click maps from different settings of the bubble radius. Plotted here are the clicks of 3 participants (b-d) who explored the same image (a) with BubbleView, but with a different bubble size: 16, 24, and 32 pixel radius, respectively. The smaller the bubble, the more clicks a participant made, and the longer the task took to complete. Overall, the same regions of interest tended to be clicked on, despite differences in bubble sizes.

participants, the similarity of BubbleView click maps to ground truth fixation maps was already 97-98% of the performance achievable in the limit. The NSS score was extrapolated to increase to 1.31 in the limit (95% C.I. [1.312, 1.315]) with a bubble size of 16, 1.32 in the limit (95% C.I. [1.320, 1.324]) with a bubble size of 24, and 1.31 in the limit (95% C.I. [1.306, 1.310]) with a bubble size of 32. As a result of these analyses, we used an average of 10–15 participants for all future BubbleView experiments.

> *Take-aways:* 10–15 participants worth of BubbleView clicks already accounted for up to 97-98% of the performance achievable in the limit of the number of participants.

*Results on ranking elements by importance*

We also explored the relationship between BubbleView clicks and eye fixations at ranking visualization elements by importance. For this purpose we used the element segmentations (e.g., title, axis, legend, etc.) available in the MASSVIS dataset [24]. For each of the 202 visualizations from Exp. 1, we overlapped the element segmentations with the fixation map of the visualization, and took the maximum value of the fixation map within the element's boundaries as its **importance score** (this analysis is inspired by the analyses used in Chapters 6.2.3-6.2.5 to evaluate the relative importance of different people, faces, and text in natural images). We averaged the element scores across all 202 visualizations to obtain an aggregate importance score for each type of element (Figure 7.7). We repeated this computation using the BubbleView click maps of the visualizations to get another set of importance scores for the same elements. The ranking of elements by importance scores according to BubbleView clicks is highly correlated to the ranking according to eye fixations (Spearman correlation = 0.96).

| Exp. 1: visualizations | Bubble Radius (pixel) | CC | NSS | Normalized NSS |
|---|---|---|---|---|
| Exp. 1.1: 51 visualizations Description task (ground-truth IOC: 1.42) | 16 | 0.86 0.87 | 1.27 1.30 | 89% ($n = 10$) 92% ($n = 38$) |
| | 24 | 0.86 0.87 | 1.27 1.30 | 89% ($n = 10$) 92% ($n = 39$) |
| | 32 | 0.86 0.87 | 1.27 1.29 | 89% ($n = 10$) 91% ($n = 40$) |
| Exp. 1.2: 51 visualizations Description task (ground-truth IOC: 1.33) | 24 | 0.82 0.84 | 1.20 1.22 | 90% ($n = 10$) 92% ($n = 20$) |
| | 32 | 0.84 0.85 | 1.20 1.22 | 90% ($n = 10$) 92% ($n = 18$) |
| | 40 | 0.83 0.84 | 1.19 1.19 | 89% ($n = 10$) 89% ($n = 11$) |
| Exp. 1.3: 100 visualizations Description task (ground-truth IOC: 1.35) | 32 | 0.84 0.84 | 1.21 1.21 | 90% ($n = 10$) 90% ($n = 10$) |

**Table 7.1.** We evaluated BubbleView clicks at approximating ground-truth eye fixations on the MASSVIS dataset by varying the bubble radius. We ran 3 sets of experiments on different subsets of the MASSVIS dataset. We measured the cross-correlation (CC) between BubbleView click maps and ground truth fixation maps, averaged over all images (CC has an upper bound of 1). The normalized scanpath saliency (NSS) score measured how well BubbleView click maps predict discrete fixation locations, averaged over all images. The NSS upper bound depends on the ground-truth data, so we included the inter-observer consistency (IOC) score of the eye tracking participants (measured using NSS). Normalizing the NSS score of the BubbleView maps by IOC allows us to report the percent of ground-truth fixations predicted by the BubbleView maps. To make the scores comparable across all the experiments, we fixed the number of participants to $n = 10$. In gray we report the results obtained by including all $n$ participants that were collected for each experiment. The difference in CC and NSS scores with different bubble radius sizes was not significant ($F < 1$ for all comparisons).

Increasing the number of BubbleView participants increases the similarity
to ground-truth eye fixations (biggest increase up to 10-15 participants)



**Figure 7.6.** The NSS score of BubbleView click maps computed with different numbers of participants, when used to predict discrete fixation locations on the MASSVIS dataset. Each point represents the score obtained at a given number of participants, averaged over 10 random splits of participants, and all 51 images used in Exp. 1.1. We include data points from 3 different bubble radius sizes. By fitting power functions of the form $an^b + c$ to each set of points, we find that these scores do not change significantly in the limit of participants ($n \to \infty$).

> *Take-aways:* BubbleView can be used to rank visualization elements by importance, predicting how often people would fixate those elements during natural viewing.

## ■ 7.3.2 Experiment 2: comparison to eye fixations on natural images

In Experiment 1 we found that BubbleView clicks offered a very good approximation to eye fixations on information visualizations with a description task. However, because free-viewing is a more common setting for human perception studies of natural images (specifically for saliency datasets), we wanted to determine if BubbleView clicks can also be used to approximate free-viewing fixations on natural images. We used similar BubbleView settings to the ones found in Exp. 1: a bubble size of 30 pixels and 15 participants worth of clicks.

> *Motivating questions*
>   • Does BubbleView generalize to natural images with a free-viewing task?

*Stimuli*
The OSIE dataset contains 700 natural images with multiple dominant objects per im-

**Figure 7.7.** (a) An example of a labeled visualization from the MASSVIS dataset. (b) By overlapping fixation maps and BubbleView click maps with such element annotations (and taking the maximum value of the map inside the element), we obtain an importance score for each element in each visualization. By averaging across 202 visualizations, we obtain an aggregate importance score per element type.

age [223]. Eye movements on this dataset were collected by instructing 15 participants to free-view each image for 3 seconds. Participants made an average of 9.3 fixations per image (3.1 fixations/sec). In this eye tracking setup, images were presented at a resolution of $800 \times 600$ pixels and 1 degree of viewing angle corresponded to 24 pixels. For our study, we randomly sampled 51 OSIE images (Figure 7.8), downsized them to $640 \times 480$ pixels, and blurred them with a sigma of 30 pixels.

*Method*
In **Exp. 2.1**, we asked participants to free-view a series of images and to click anywhere they want to look for 10 sec per image. We used a bubble radius of 30 pixels, equal to about 1.5 degrees of visual angle in the eye tracking study. Although the viewing time for the OSIE eye tracking study was 3 sec per image, we increased this time for the BubbleView experiment to account for the time of clicking a mouse. We piloted different viewing times and determined 10 sec to be appropriate (clicking took about 3 times as long as natural viewing). We collected an average of 60 participants worth of BubbleView click data for each image.

Apart from ground truth eye fixations, mouse movements using the related SALI-CON methodology are also available for the OSIE dataset [103]. To facilitate a direct comparison between BubbleView and SALICON, in **Exp. 2.2** we re-ran data collection with BubbleView, replacing mouse clicks with mouse movements, with a bubble radius of 30 pixels. As in SALICON, we used a task time of 5 seconds. The results of this experiment are discussed in Section 7.4.2, in the context of other comparisons to the

**Figure 7.8.** Example images from the OSIE dataset. Dataset images (a), with corresponding ground-truth fixation maps (b) and BubbleView click maps (c). We show cases where BubbleView maps have high similarity, and cases with low similarity, to fixation maps.

SALICON methodology.

*Results*

During 10 seconds of viewing, participants made an average of 13.1 clicks, or about 1.3 clicks/sec - three times fewer clicks than fixations per second.

In Exp. 2.1, the similarity between BubbleView click maps and ground truth fixation maps with free-viewing on natural images was smaller (NSS = 2.61, CC = 0.81, Table 7.2) than in Exp. 1 with visualizations. Even though eye tracking participants are quite consistent with each other on the OSIE dataset (IOC = 3.35), BubbleView participants are not as predictive of eye tracking participants in this case. BubbleView clicks of 54 participants can predict 80% of eye fixations, while the projected performance in the limit only converges to 82% (95% C.I. [2.742, 2.754]). However, 10 BubbleView participants can already account for 78% of eye fixations.

Exp. 2.2 showed that a related methodology using a moving-window approach [103] is no better at approximating ground-truth eye fixations on this dataset (Table 7.2). In fact, to achieve the same performance as BubbleView, SALICON actually requires more participants (Section 7.4.2). BubbleView can serve as an affordable and scalable alternative. When running a large number of eye tracking experiments is infeasible, BubbleView can be used for studying human perception and collecting large-scale saliency datasets (e.g., for training models as in Chapter 9).

*Take-aways:* Similarity between BubbleView clicks and eye fixations is lower on natural images with a free-viewing task than with visualizations with a description task. Despite this, 10 BubbleView participants can already account for 78% of eye fixations on natural images, so BubbleView can still serve as an affordable approximation to eye tracking.

**Table 7.2.** We evaluated BubbleView clicks at approximating ground-truth eye fixations on the OSIE dataset. We ran BubbleView data collection using mouse clicks (Exp. 2.1) and using mouse movements (Exp. 2.2). For comparison, we also include the performance of the SALICON methodology on the same dataset (Section 7.4.2). For fair comparison with in-lab SALICON, we only used $n = 12$ participants per image per study. The difference in scores at $n = 12$ participants was not significant [F(200)=1.81, n.s.]. In gray we report the results obtained by including all $n$ participants that were collected for each experiment.

| Exp. 5.3: natural scenes (ground-truth IOC: 3.35) | CC | NSS | Normalized NSS |
|---|---|---|---|
| BubbleView (clicks) | 0.81 | 2.61 | 78% ($n = 12$) |
| | 0.84 | 2.69 | 80% ($n = 54$) |
| BubbleView (movements) | 0.81 | 2.52 | 75% ($n = 12$) |
| | 0.83 | 2.55 | 76% ($n = 49$) |
| SALICON | 0.81 | 2.52 | 75% ($n = 12$) |
| | 0.84 | 2.61 | 78% ($n = 92$) |
| In-lab SALICON | 0.81 | 2.61 | 78% ($n = 12$) |
| | 0.81 | 2.61 | 78% ($n = 12$) |

### ■ 7.3.3 Experiment 3: comparison to eye fixations on static webpages

Apart from natural images, webpages are another image type that frequently serve as the focus of eye tracking and usability studies [29, 39, 151, 184, 194, 195]. For this reason, we wanted to test the generalizability of the BubbleView methodology to webpages. Because the static webpage images were denser in visual and information content than the information visualizations and natural images from the first two experiments, we evaluated a number of different BubbleView settings to try to find the best approximation to eye fixations. We varied bubble radius size and viewing time. As in the original FiWI eye-tracking experiment, we started with a free-viewing task. Similar to Exp. 1, we also tried a description task with unlimited task time.

**Figure 7.9.** Example images from the FiWI dataset. Dataset images (a), with corresponding ground-truth fixation maps (b) and BubbleView click maps (c). We show cases where BubbleView maps have high similarity, and cases with low similarity, to fixation maps.

*Motivating questions*

- Does BubbleView generalize to webpages?

- How do the task and viewing time affect performance?

- Does viewing time interact with bubble size?

*Stimuli*

The FiWI dataset contains 149 screenshots of static webpages collected from various sources on the Internet and sorted into pictorial (dominated by pictures such as photo sharing websites), text (high density text such as encyclopedia websites), and mixed types [194]. Eye movements on this dataset were collected by instructing 11 participants to free-view each webpage for 5 seconds. Participants made an average of 17.9 fixations per image (3.6 fixations/sec). In this eye tracking setup, 1 degree of visual angle was approximately 50 pixels.

We sampled 17 images from each of the three categories (pictorial, text, mixed), resulting in a total of 51 images (Figure 7.9). We downsized the images from $1360 \times 768$ pixels to $1000 \times 565$ pixels to fit within a typical MTurk browser window, while preserving image aspect ratios. These webpages tended to have more varied font size compared to the images in Exp. 1–2. We manually selected a blur sigma of 50 pixels to distort the text on these images beyond legibility.

*Method*

We ran experiments with two task types where participants were asked to either free-view or describe each webpage. In **Exp. 3.1**, with the free-viewing task, we used a 2 x 3 factorial design (viewing time: 10 sec or 30 sec; bubble radius: 30, 50, or 70 pixels). In **Exp. 3.2**, with the description task, we used a bubble radius of 30 pixels

Longer task time increases the similarity between BubbleView clicks
and eye fixations. A defined description task outperforms free-viewing.



**Figure 7.10.** The NSS score of BubbleView click maps computed with different numbers of participants, when used to predict discrete fixation locations on the FiWI dataset. Each point represents the score obtained at a given number of participants, averaged over 10 random splits of participants, and all 51 images.

and unlimited time. We collected an average of 15 participants worth of BubbleView click data for each image under each task.

*Results on stimuli*
In the free-viewing task (Exp. 3.1), participants made an average of 1.0–1.8 clicks/sec, while in the description task (Exp. 3.2), participants made an average of 0.5 clicks/sec, indicating that they spent more than half the time typing descriptions. Clicking took about 3 times longer than natural viewing. As in Exp. 1, the number of clicks per second monotonically decreased with increasing bubble size, even though viewing time was fixed (Exp. 3.1). Tripling the viewing time from 10 to 30 seconds did not quite triple the number of clicks, but increased them by 2.2–2.6 times.

The similarity between BubbleView click maps and ground truth fixation maps on webpages was lowest of all image types tested so far in Exp. 1–3 (Table 7.3). However, the inter-observer consistency of eye tracking participants is also lowest on webpages (IOC = 1.85). Recall that IOC between eye tracking participants serves as an upper bound for how well BubbleView clicks can predict eye fixations. After accounting for IOC, the normalized NSS scores show that BubbleView clicks can account for up to 78% of eye fixations on webpages, similar to the score on natural images (Exp. 2).

IOC was highest on the all-text webpages (NSS = 1.97), followed by the pictorial (NSS = 1.77) and mixed (NSS = 1.80) webpages. While the difference in NSS scores was not significant across webpage types for the similarity between BubbleView clicks and eye fixations, the NSS scores were consistently higher for the text webpages. Only for one case, with a bubble size of 30 pixels and 10 seconds of viewing, were the NSS scores for the pictorial webpages the highest. This provides evidence that clicks tend to be more consistent with fixations on text elements.

> *Take-aways:* Both fixation and click data is more varied on webpages. Webpage images with lower IOC scores (more eye tracking variability) also had worse BubbleView similarity scores. Normalizing for IOC, BubbleView clicks can account for 78% of eye fixations on webpages (as for natural images).

*Results on time, bubble size, and task*
We ran a two-way ANOVA (time $\times$ bubble size) on Exp. 3.1. The main effect of time on the CC and NSS scores was significant [CC: $F(1,300)=19.25$, $p < .01$, NSS: $F(1,300)=9.65$, $p < .01$], respectively but the effect of bubble size was not [CC: $F(2,300)=1.92$, n.s., NSS: $F(2,300)=1.14$, n.s.]. The interaction effect between time and bubble size was significant [CC: $F(2,300)=6.95$, $p < .01$, NSS: $F(2,300)=3.35$, $p < .05$].

With a viewing time of 10 seconds, a bubble size of 30 pixels was too small, achieving significantly lower CC scores than bubble sizes of 50-70 pixels ($p < .05$). With a viewing time of 30 seconds, however, a bubble size of 70 pixels was too large, achieving significantly lower CC scores than bubble sizes 30-50 pixels ($p < .01$). No significant differences were found among the NSS scores (Table 7.3). There exists a trade-off: with a longer viewing time, a smaller bubble radius provides more consistent clicks among participants; when limited by a shorter time, a larger bubble size becomes necessary.

Given a bubble size of 30-50 pixels, the CC scores were significantly higher for a task duration of 30 seconds compared to 10 seconds ($p < .05$). The difference in NSS scores was only significant for the bubble size of 30 pixels. No significant differences were found with a bubble size of 70 pixels. Overall, BubbleView click maps generated with longer task durations of 30 seconds or longer (including with a description task) better approximated eye fixations than with a 10 second task duration. From this we conclude that information-dense images like websites require either longer viewing times or better defined tasks than free-viewing.

From Exp. 3.2, we found that for small numbers of participants ($n < 12$), the description task generated BubbleView click maps more similar to ground-truth eye fixations than the free-viewing task under all settings (Figure 7.10). The difference between the tasks is larger for smaller number of participants, and decreases with each extra participant. The click data tends to converge faster when a targeted task like description is used. However, this advantage disappears with more participants and a longer task time (30 sec, 30 pixel bubble radius). A description task takes longer and is more expensive to run, but might be a better choice when few participants are available.

**Table 7.3.** We evaluated BubbleView clicks at approximating ground-truth eye fixations on the FiWI dataset. BubbleView maps were computed with 12 participants for all experiments below. The score of the BubbleView maps predicting the ground-truth fixation maps is reported in CC, and the score of the BubbleView maps predicting the discrete fixation locations is reported in NSS. Normalized NSS is calculated by normalizing the NSS score by the inter-observer consistency (IOC) of the eye tracking participants.

| Exp. 3: webpages (ground-truth IOC: 1.85) | Time (sec) | Bubble Radius (pixel) | CC | NSS | Normalized NSS |
|---|---|---|---|---|---|
| Free-viewing | 10 | 30 | 0.52 | 1.20 | 65% |
| Free-viewing | 10 | 50 | 0.57 | 1.34 | 72% |
| Free-viewing | 10 | 70 | 0.56 | 1.30 | 70% |
| Free-viewing | 30 | 30 | 0.63 | 1.45 | 78% |
| Free-viewing | 30 | 50 | 0.61 | 1.41 | 76% |
| Free-viewing | 30 | 70 | 0.57 | 1.32 | 71% |
| Description | unlim. | 30 | 0.63 | 1.46 | 79% |

*Take-aways:* the less viewing time available, the larger the bubble size should be in order to better approximate free-viewing fixations. For a study with fewer participants, a description task is better than a free-viewing task.

## ■ 7.4  Experiments comparing BubbleView to related methodologies

## ■ 7.4.1  Experiment 4: comparison to importance annotations on graphic designs

We hypothesized that the regions on an image where participants click using the BubbleView methodology correspond to the most important regions of the image. To test this hypothesis, we used the GDI dataset [154] which comes with explicit importance annotations, where participants were instructed to annotate the image regions they considered important in graphic designs. We used this dataset to evaluate whether the number of BubbleView clicks on image regions corresponds to explicit judgements of importance.

*Motivating questions*
- Does BubbleView generalize to graphic designs?

- Do BubbleView clicks correspond to regions of importance on graphic designs?

**Figure 7.11.** Example images from the GDI dataset. Images from the dataset (a), along with the provided explicit importance annotations (b). We show cases where BubbleView maps have high correlation, and cases with low correlation, to the importance annotations, in terms of how design elements are ordered by importance (c).

*Stimuli*

The Graphic Design Importance (GDI) dataset contains 1,075 single-page graphic designs (e.g., advertisements, flyers, and posters consisting of text and graphical elements), collected from Flickr [154]. No eye movements were collected for this dataset. O'Donovan et al. [154] highlighted two downsides of eye movements for this type of data: (1) fixations vary significantly over individual elements (like text blocks) even though those regions should have a uniform importance, and (2) eye fixations may occur in unimportant regions as a design is scanned and do not reflect conscious decisions of importance. Instead, 35 MTurk participants were asked to label important regions with binary masks, and these masks were averaged over all participants to produce a final importance map per design. O'Donovan et al. [154] noted that although importance maps produced by individual users are noisy, the average map gives a plausible relative ranking over design elements.

We sampled 51 images from the GDI dataset at the original resolution of $600 \times 400$ pixels (Figure 7.11). We blurred the images with a sigma of 30 pixels, manually chosen to distort text beyond recognition.

*Method*

We ran an experiment with a bubble radius of 50 pixels and viewing time of 10 seconds, in which participants were asked to free-view each graphic design. BubbleView strikes

a balance between eye fixations and explicit importance judgements for these images: (1) like fixations, clicks are collected in a free-viewing setting and are not uniform over design elements, but (2) like explicit annotations, the decisions of where to click reflect conscious decisions of importance. We collected an average of 15 participants worth of BubbleView click data for each image.

*Analysis*

Unlike the quantitative evaluations in the previous sections, we did not directly compare the BubbleView click maps to the graphic design importance (GDI) maps. The spatial distributions of the explicit importance annotations in the GDI dataset are different from the click maps generated by our methodology. By construction, the importance annotations are uniform over design elements in the GDI dataset, while BubbleView clicks are not. For a fairer comparison, we computed the importance values each methodology assigns to different elements within each design (similar to the analysis at the end of Section 7.3.1 used to rank visualization elements by importance).

We used bounding boxes to manually annotate all the elements in the 51 graphic designs chosen. For each design we normalized the GDI ground-truth importance map and the BubbleView click map. We took the maximum value of each map within an element's bounding box as the importance score of that element. We correlated the importance scores assigned by both methodologies to the elements in each design (Figure 7.12).



**Figure 7.12.**   Importance maps were overlapped with element bounding boxes (outlined in blue) and the maximum map value per box was taken to be the importance score for that element (scores are the numbers above each box). Maps were first normalized to have values between 0 and 1, so the importance scores for all the graphic design elements also fall within the same range, where 1 corresponds to the most important element. In the case of the GDI importance map, MTurk workers made explicit judgements about aspects of the graphic design they considered the most important. A region of a graphic design has an importance score of 1 if all MTurk workers labeled that element as important. In the BubbleView study, MTurk workers clicked a blurred graphic design to expose small regions of the design at full resolution. A region of a graphic design has an importance score of 1 if the density of MTurk clicks in that region was highest.

*Results*

Across all 51 graphic designs, we achieved an average Pearson correlation of 0.66 and an average Spearman (rank) correlation of 0.60 between the element importance scores as assigned by BubbleView versus the original GDI annotations. Over 70% of graphic designs had a correlation over 0.4. BubbleView importance maps can reasonably approximate explicit importance judgements for ranking elements of graphic designs, although there are some differences. For instance, the blurring of the image may interfere with visual features seen at different scales, as in the last two example images in Figure 7.11. Depending on the blur, certain visual elements might not be clicked on (e.g. in Figure 7.11, the *note* because it blended into the background when blurred; the *eye* because it was already visible in the blurred version).

> *Take-aways:* BubbleView can be used to rank graphic design elements by importance. However, due to the varied feature sizes, blurring might significantly impact which design regions are clicked.

## ■ 7.4.2 Experiment 5: comparison to mouse movements on natural images

The most similar methodology to BubbleView is SALICON [103], which was introduced at roughly the same time[5]. SALICON is also intended to be used in a crowdsourcing setting to approximate eye fixations [103]. The differences are that SALICON captures continuous mouse movements, instead of clicks, and images are blurred adaptively, with a multi-resolution blur, recomputed for each cursor position. We investigated whether BubbleView click maps are similar to SALICON mouse movement maps, when averaged over multiple participants. Because the SALICON blur is multi-resolution and adaptive, we experimented with different blur sigmas and bubble sizes in BubbleView, to find a fixed setting of parameters that best approximates the SALICON viewing conditions. We also compared SALICON and BubbleView at approximating eye fixations collected in a controlled lab setting, since both methodologies are presented as alternatives to eye tracking.

> *Motivating questions*
> - Under what settings does BubbleView most closely match SALICON?
> - Which methodology better approximates eye fixations on natural images?

*Stimuli*

The SALICON dataset consists of mouse movements collected on 20K MS COCO (Microsoft Common Objects in Context) natural images [138]. In the original study, mouse

---

[5]The SALICON and BubbleView methodologies were introduced a few months apart, but to different communities: Jiang et al. [103] to computer vision and Kim et al. [112] to human-computer interaction.

**Figure 7.13.** Example images from the SALICON dataset. Example dataset images (a), and ground truth mouse movements collected by SALICON (b). We show cases where BubbleView maps have high similarity, and cases with low similarity, to SALICON maps (c).

movements were collected on Amazon's Mechanical Turk by presenting images to participants for 5 seconds each and allowing them to freely explore each image by moving the mouse cursor. We randomly sampled 51 images at the original image size of 640 × 480 pixels from the SALICON dataset (Figure 7.13).

*Method*

In **Exp.  5.1**, we used a 3 × 3 factorial design (blur sigma: 30, 50, and 70 pixels; bubble radius: 30, 50, and 70 pixels; see Figure 7.14). Using a free-viewing task, we had participants explore each image for 10 seconds each. We wanted to account for longer times to click, rather than move, the mouse.

To disentangle the influence of mouse clicks/movements versus fixed/adaptive blur on the methodology differences between SALICON and BubbleView, we ran **Exp. 5.2**, using BubbleView with a moving-window approach like SALICON, but maintaining a fixed blur kernel. In this setup participants used mouse movements to reveal image regions at normal resolution. We had two experiment conditions (bubble radius sizes of 30 and 50 pixels) with a fixed blur sigma of 30 pixels (found appropriate in Exp. 5.1) and viewing time of 5 seconds (as in SALICON). We collected an average of 15 participants worth of BubbleView click data for each image under each condition.

*Results on using BubbleView to approximate SALICON*

We ran a two-way ANOVA (blur × bubble size) on Exp. 5.1. The main effect of bubble size was not significant [CC: $F(2,450)=2.28$, NSS: $F(2,450)=0.19$, n.s.] (as found in Exp. 1 and 3.1). The main effect of blur on scores was significant [CC: $F(2, 450)=19.97$, $p < .01$, NSS: $F(2,450)=6.86$, $p < .05$]. BubbleView with a blur radius of 70 pixels

(a) BubbleView
with fixed blur

(b) SALICON
with adaptive multi-resolution blur

**Figure 7.14.** We used 9 different parameter settings in our BubbleView experiments, on images from the SALICON dataset (a). We wanted to find a fixed setting of bubble size and blur to mimic the adaptive multi-resolution blur used in the SALICON methodology (b). The rightmost figure is from Jiang et al. [103]; ©Martijn van Exel.

**Table 7.4.** We evaluated BubbleView click maps (with $n = 12$ participants per image) at approximating SALICON mouse movements, measured using CC and NSS metrics. Normalized NSS is computed by taking into account the IOC of the SALICON participants (NSS = 1.50). Both bubble radius and blur sigma are measured in pixels. BubbleView with a blur radius of 70 pixels achieved significantly lower CC scores than with other blur settings ($p < .01$ for all bubble sizes). The other differences were not significant.

| | | | Blur Sigma (pixel) | | |
| --- | --- | --- | --- | --- | --- |
| | | | 30 | 50 | 70 |
| | | CC | 0.84 | 0.84 | 0.78 |
| | 30 | NSS | 1.21 | 1.15 | 1.06 |
| | | Normalized NSS | 81% | 77% | 71% |
| Bubble radius (pixel) | | CC | 0.86 | 0.84 | 0.80 |
| | 50 | NSS | 1.23 | 1.15 | 1.04 |
| | | Normalized NSS | 82% | 77% | 69% |
| | | CC | 0.84 | 0.84 | 0.79 |
| | 70 | NSS | 1.20 | 1.11 | 1.04 |
| | | Normalized NSS | 80% | 74% | 69% |

achieved significantly lower CC scores than with other blur settings ($p < .01$ for all bubble sizes). We did not find an interaction effect between blur and bubble size [CC: $F(4,450)=0.44$, NSS: $F(4,450)=0.04$, n.s.]. We found highest similarity between BubbleView click maps and SALICON maps at bubble radius sizes of 30–50 pixels and blur sigma of 30–50 pixels (Table 7.4), for which the normalized NSS scores ranged from 77% to 82%.

What are the remaining differences? Using mouse movements, more points of interest are generated than using clicks. Many of the points sampled using mouse movements occur in the transition between regions in an image, and might be introducing noise into the data (Figure 7.15). This suggests that a different threshold might be more effective at converting continuous mouse movements into discrete points of interest. An advantage of the BubbleView clicks is that no such post-processing is necessary, since the clicks directly correspond to points of interest.

In Exp. 5.2, we modified BubbleView to collect continuous mouse movements and shortened the time per image to 5 sec, such that the only remaining difference with SALICON was the treatment of blur. We observed that the mean number of samples was 143.02 (SD=13.14) using the sampling rate of 100 Hz, which translates to 14,302 raw samples, on average, per participant. This is significantly larger than the mean click count of 13.09 (SD=1.38) per participant in Exp 5.1.

With the moving-window BubbleView setting the scores were: for bubble size 30: CC: 0.87, NSS: 1.21, normalized NSS: 81%; bubble size 50: CC: 0.88, NSS: 1.24, normalized NSS: 83%. Compared to the clicks, these scores were not statistically significantly different [$F(200) < 2.2$, n.s.]. In other words, BubbleView can approximate SALICON with or without mouse movements. Importantly, BubbleView can approximate SALICON without requiring a multi-resolution adaptive blur, simply with a single fixed blur setting. Our fixed blur setting is much less computationally expensive and does not require the pre-study system checks as in Jiang et al. [103].

> *Take-aways:* BubbleView with a bubble size of 30–50 pixels and a blur sigma of 30–50 pixels can approximate the continuous mouse movements and adaptive, multi-resolution blur of the SALICON methodology.

*Results on using both methodologies to approximate eye fixations*
In Exp. 2.2 we compared BubbleView clicks and mouse movements to SALICON mouse movements at approximating ground truth eye fixations on 51 OSIE images. The BubbleView click maps (with $n = 12$ participants, bubble radius of 30 pixels) achieve NSS $= 2.61$ (CC $= 0.81$) at predicting ground-truth fixation maps, compared to SALICON mouse movement maps which achieve NSS $= 2.52$ (CC $= 0.81$). It takes over 30 SALICON participants to achieve the same similarity to fixation maps as 12 BubbleView participants (Figure 7.16). Replacing BubbleView clicks with mouse movements actually decreases performance: NSS $= 2.52$ (CC $= 0.81$), but this drop in performance is not significant at the $p = .05$ level. For all feasible numbers of participants ($n < 60$ in Figure 7.16), BubbleView offers a better approximation to eye fixations than SALICON.

(a) points of interest from SALICON mouse movements



(b) points of interest from BubbleView mouse clicks



**Figure 7.15.** When participants can move the mouse anywhere on the image without having to click, the collected data contains motion traces as byproducts (a). Instead of only capturing the points of interest in an image where an observer's attention stops, the moving-window approach also captures the transitions between these regions, which are less relevant and add noise to the data. Although these trajectories can be post-processed into discrete regions of interest, our approach is to directly collect participant mouse clicks on points of interest, with no further post-processing required (b).

Data was also available for 12 in-lab participants who used the SALICON methodology to view images in a controlled lab setting [103]. The in-lab SALICON maps, which capture these mouse movements, achieve NSS = 2.61 (CC = 0.81) when compared to fixation maps, the same score as our BubbleView maps (Table 7.2). From Figure 7.16 we can see that the performance of the in-lab SALICON is increasing at a greater rate than either BubbleView or online SALICON. However, more in-lab SALICON participants would be needed to see whether this trend continues. In any case, it requires a controlled lab setting, which we aim to avoid.

*Take-aways:* On a natural image dataset, BubbleView clicks better approximate eye fixations than SALICON mouse movements for all feasible numbers of participants ($n < 60$). BubbleView performed better with clicks than BubbleView with mouse movements.

**Figure 7.16.** The NSS score obtained by comparing mouse clicks and mouse movements to ground truth eye fixations on natural images in the OSIE dataset. We compare mouse clicks gathered using BubbleView on MTurk (purple), mouse movements gathered using BubbleView on MTurk (green), mouse movements gathered using SALICON on MTurk (blue), and mouse movements gathered using SALICON in a controlled lab setting (black crosses). Each point represents the score obtained at a given number of participants, averaged over 10 random splits of participants and all 51 images used.

## ■ 7.5 Addendum: free-viewing versus description

After publication of [114], we received questions about whether the description task used with information visualizations was essential to the results we were seeing in Experiment 1 (Sec. 7.3.1) - specifically, that 90% of free-viewing fixations could already be accounted for by the BubbleView clicks of 10 participants (who were required to provide a description as they clicked around). In this addendum to the paper we address this question:

> *Motivating question*
> - How does the task of free-viewing versus description affect the clicks of BubbleView participants?

*Stimuli*
We sampled 30 images from the set used in Exp. 1.3 (Sec. 7.3.1). Images were a maximum dimension of 600 pixels to a side and were blurred with a sigma of 40 pixels. We used a bubble radius of 32 pixels.

*Method*
In a single HIT, participants were shown all 30 images, for 10 seconds each and 2 sec-

**Table 7.5.** We evaluated BubbleView clicks at approximating ground-truth eye fixations on 30 images from the MASSVIS dataset, using 2 tasks: description and free-viewing. BubbleView maps were computed with $n = 12$ participants in black, and including all $n$ participants collected for each experiment in gray. The score of the BubbleView maps predicting the ground-truth fixation maps is reported in CC, and the score of the BubbleView maps predicting the discrete fixation locations is reported in NSS. Normalized NSS is calculated by normalizing the NSS score by the inter-observer consistency (IOC) of the eye tracking participants on the 30 images.

| Visualizations (ground-truth IOC: 1.42) | CC | NSS | Normalized NSS |
|:---:|:---:|:---:|:---:|
| Description task | 0.87 | 1.30 | 92% ($n = 12$) |
| Free-viewing task | 0.72 | 1.07 | 75% ($n = 12$) |
| | 0.75 | 1.13 | 80% ($n = 57$) |

onds between consecutive images, with the instructions to "click anywhere you want to look". The resulting clicks collected correspond to a free-viewing task. We collected an average of 60 participants worth of BubbleView click data for each image.

*Results*

In the description task (from Sec. 7.3.1), participants made an average of 64 clicks per image over a viewing interval of approximately 3 minutes/image. If we account for the time to write descriptions, participants spent 15-30% of the total task time (of an average of 3 minutes/image) clicking, or approximately 27-54 seconds. That is, 1-2.4 clicks/sec. In comparison, in the free-viewing task, participants made an average of 15 clicks per image over a 10 second viewing interval, or 1.5 clicks/sec. In the original free-viewing task with an eyetracker, participants made an average of 40 fixations per image in 10 seconds of viewing, or 4 fixations/sec, attending to 2-3 times more locations than with clicking.

From Table 7.5 we see that for the same number of participants, the click data obtained under a description task is significantly more similar to eye fixations (under free-viewing) than the click data obtained under a free-viewing task. With 12 participants, the click maps obtained from the free-viewing task achieve an NSS score of 1.07, extrapolated to increase to 1.19 (95% C.I. [1.18, 1.19]) in the limit (Fig. 7.17), compared to an NSS of 1.30 already achievable with 12 participants performing a description task (and a limiting NSS score of 1.37, 95% C.I. [1.35, 1.38]). In other words, increasing the number of participants can not compensate for the difference in the quality of clicks generated between the two tasks. Given a description task, participants explore more of the visualization than with the free-viewing task (where participants focused predominantly on titles and main text elements). Note that this result can be confounded with the amount of time participants were given to freely explore each infographic (limited to 10 seconds) in the free-viewing scenario. However, without a task, a longer task may not necessarily keep online crowdworkers motivated.

BubbleView clicks from a description task are more similar to
free-viewing eye fixations than clicks from a free-viewing task



**Figure 7.17.** The NSS score obtained by comparing mouse clicks under free-viewing and description tasks to ground truth eye fixations on a small set of visualizations in the MASSVIS dataset. Each point represents the score obtained at a given number of participants, averaged over 10 random splits of participants and all 30 images used.

> *Take-aways:* A description task is recommended over a free-viewing task for BubbleView clicks to better approximate (free-viewing) eye fixations collected using an eye tracker.

## ■ 7.6  Summary of experimental results

**Similarity of BubbleView clicks to eye fixations:** We showed that across 3 different image types (information visualizations, natural images, and static webpages) and 2 types of tasks (free-viewing and description), BubbleView clicks provide a reasonable approximation to eye fixations collected in a controlled lab setting. Specifically, across all these image types BubbleView clicks accounted for over 75% of eye fixations when only 10–15 BubbleView participants were used (Tables II-IV). Of all settings, BubbleView clicks provided the best approximation to eye fixations on information visualizations with a description task, accounting for up to 90% of eye fixations with only 10 participants, and 92% with 20 participants (Table II). On both natural images and websites, BubbleView clicks could account for up to 78% of eye fixations with 10–12 participants (Tables III, IV). The fixations of eye tracking participants were much more consistent on the natural images than on the websites, so the viewing behavior on natural images should be easier to predict. Despite the remaining gap between BubbleView

**Figure 7.18.** Taking a horizontal cross-section of the average BubbleView click map and the average fixation map across 51 images on 3 datasets, we see the fixation map has a consistent center bias. This replicates the analysis used by Tatler [206] to report on human fixation bias in natural images. This bias emerges as a peak near the center of an image, which corresponds to the midway point along the $x$-axis in each of these plots. The BubbleView click map does not have this bias, which accounts for some of the systematic differences observed between the click and fixation maps. At the same time, the bubble clicks tend to capture the same general characteristics as fixations, for instance of increased attention in the leftmost parts of visualizations and webpages, corresponding to the titles and headers.

clicks and eye fixations for natural images and webpages, the fact that already 10–15 BubbleView participants achieves a reasonable approximation to fixations is promising for perception studies that might otherwise require specialized eye-tracking hardware.

**Remaining differences between BubbleView clicks and eye fixations:** Part of the remaining gap between BubbleView clicks and eye fixations is that BubbleView does not capture the unconscious movements of the eyes due to bottom-up, pop-out effects, or systematic biases. One such systematic bias commonly referred to in the eye tracking literature is center bias [18, 32, 206], whereby a relatively high number of fixations occur near the center of the image. One explanation for such bias is that it is part of an optimal viewing strategy that is involved in planning successive fixations. By averaging fixation maps across dataset images, we can see a peak near the spatial center of the image emerge across the eye fixations, but not the BubbleView clicks (Figure 7.18). Because BubbleView naturally slows down the exploration task by making participants consciously decide where to click next, it captures higher-level viewing behaviors not as affected by systematic biases. We recommend using BubbleView with a well-defined task, like describing the content of the visual input, to measure which regions of that visual input are most important or relevant for the task.

A recent paper by Tavakoli et al. [209] analyzes some of the semantic differences between eye fixations and mouse movements on the OSIE dataset, by taking into account annotated image regions. They find that there tends to be more disagreement between eye fixations and mouse movements in background regions of the image.

**Effect of BubbleView parameters:**  The BubbleView click maps were quite robust across different parameter settings. We did not find significant effects of bubble radius on the resulting BubbleView clicks (Exp. 1,3,5). Across all our experiments (Exp. 1–5), we found that a blur kernel sigma in the range of 30–50 pixels was appropriate for all of our image types, where we manually selected a sigma value for each image dataset to ensure that text was unintelligible when blurred and would require explicit clicking on to read. In other words, to mimic peripheral vision, the blur level was chosen to eliminate legible details beyond the focal region. However, a blur sigma with a 70-pixel radius was too high, and seemed to hinder exploration of the image by eliminating too much context, as similarity of BubbleView clicks to eye fixations significantly dropped for this blur level compared to a blur of 30–50 pixels (Exp. 5).

We found that a bubble radius in the range of 30 to 50 pixels seems to consistently work best for different image types and image sizes that comfortably fit within the browser window (ranging from $500 \times 500$ to $1000 \times 600$ pixels). Here "best" refers to the ability of BubbleView clicks to most closely approximate fixations on images with the smallest number of participants. Smaller bubble sizes lengthened the duration and effort for completing the task, for the same quantitative results. Our chosen bubble sizes typically corresponded to 1–2 degrees of visual angle as measured in the corresponding eye tracking experiments. A bubble size of 1–2 degrees of visual angle mimics the size of the foveal region during natural viewing.

However, bubble radius is also intricately related to task timing and image complexity (Exp. 3). The more content there is on an image to look at, the more time that is required; the smaller the bubble, the more clicks to explore all of the content. A larger bubble radius can compensate for less available time, because each click exposes more of the image. For best results, we recommend a smaller bubble radius but longer task time. In our studies, the longest time for free-viewing tasks was 30 seconds (Exp. 3). For description tasks, participants spent an average of 1.5–3 minutes per image, clicking and describing (Exp. 1,3).

The number of clicks participants made decreased with increasing bubble size, even though the time for the task stayed the same. We observed this trend across all of our experiments. On average, 1–1.5 clicks were made per second in the BubbleView setup, compared to an average of 2–3 fixations per second in eye tracking studies. The BubbleView setup (when implemented with clicks) slows down visual processing so about half as many interest points are examined every second.

The best prediction performance overall occurs in the setting of a well-defined task, such as describing the visual content of an image. However, tasks must be well-matched to the images used. For instance, asking participants to describe an information visualization is well-defined because each of the visualizations we used had a main message that was being communicated (Exp. 1). On the other hand, we did not use the description task for the graphic designs (Exp. 4), because it was harder to objectively define what should be described.

**Number of participants, task, and data quality:** For our tasks, we found 10–15 participants sufficient, accounting for over 97% of the performance achievable with 40 participants (Exp. 1,2), where performance is measured by how many of the eye fixations can be approximated by clicks. The more participants, the better the data, as noisy clicks get averaged out. However, when there is a constraint on how many participants can be recruited/afforded, a more involved task (like asking the participant to provide a text description) can result in cleaner data (Exp. 3). Such a task adds an energy barrier to clicking: to minimize effort, participants are more likely to click on image regions informative for completing the task, rather than randomly.

**Mouse clicks versus movements:**  We compared our methodology of collecting discrete mouse clicks to SALICON's moving-window approach [103] in Exp. 5. We found that for any number of participants less than 60, BubbleView is a better approximation to ground truth eye fixations (Figure 7.16). This is similar to the task, data-quality trade-off discussed above. Clicks add an energy barrier to action: since clicking takes more effort than moving the mouse, participants are more selective about where they click. As a result, BubbleView provides cleaner data with fewer artifacts, such as the byproducts of continuous mouse movements (Figure 7.15). Furthermore, the moving-window methodology requires post-processing to differentiate mouse positions corresponding to points of interest from transitions. Collecting clicks directly eliminates such post-processing steps.

On the other hand, a byproduct of the higher effort of clicking on an image area rather than moving a mouse over it, is that fewer image areas will be explored by clicking. If the focus of the study is to select the most important regions in an image, then clicks should suffice. In Table 7.6 we summarize the tradeoffs between the two methodologies. We note additionally that we were able to approximate SALICON's multi-resolution adaptive blur with a single, fixed blur (Exp. 2,5) to achieve similar performances at much lower computational cost.

**Table 7.6.** Comparison of BubbleView and SALICON [103]. SALICON consists of capturing continuous mouse movements on an image with adaptive multi-resolution blur. The blur is continuously recomputed for every mouse location at 100 Hz. Continuous mouse tracks are discretized into points of interest using experimenter-specified thresholds. In BubbleView, discrete mouse clicks are collected on an image with a fixed blur. This is easier to implement and has fewer computational limitations. No additional post-processing is required. The collected BubbleView data is less noisy and converges faster, although clicking takes more time.

| Property | BubbleView | SALICON |
| --- | :---: | :---: |
| Speed of convergence to eye fixations | faster | slower |
| Number of participants required | fewer | more |
| Time per task | higher | lower |
| Post-processing | less | more |
| Computational cost | less | more |

**BubbleView for image importance:** The density of clicks in different image regions roughly corresponds to the importance of those regions. Specifically, across a collection of graphic designs, BubbleView clicks on different design elements correlated with explicit importance judgements made on the same designs (Exp. 4). BubbleView clicks ranked visualization elements similarly to eye fixations (Exp. 1). Thus, BubbleView can be used not only to derive conclusions about human perception (where people look), but also to make general conclusions about images and designs: how is importance distributed across an image? Which design elements are most important? This knowledge can in turn can be leveraged for design applications (Chapter 10).

**Data quality and filtering:** BubbleView participants were quite consistent with each other in where they clicked, leading to a relatively fast convergence of the aggregate BubbleView click maps to ground truth eye fixation maps. For most of our experiments, we found about 10-15 participants provided enough click data to reasonably approximate eye fixations.

After collecting the BubbleView data, we performed a number of filtering steps, including throwing out participants who did not click a minimum number of times and additional clicking outliers. This filtering of participants and bubbles lead to a data reduction of only 2% on average, indicating that initial data quality was pretty high.

The description task has the additional benefit of providing another filtering layer: if a participant-provided description is evaluated as poor, we can assume that they did not do the task with sufficient thoroughness, or clicked in regions of the image that were irrelevant for the task. This filtering step can either be performed manually by the experimenter or implemented as a crowdsourcing task (e.g., by having Amazon Mechanical Turk workers rate descriptions by quality).

**Cost:** The price to obtain a BubbleView click map per image depends on the amount of time a participant spends on each image and the total number of participants recruited. The average hourly rate for Amazon's Mechanical Turk is \$6/hour, so we use \$0.1/min for our tasks. It is common to make MTurk tasks bite-sized (e.g., a few minutes to 10–15 min each) [116]. Using these guidelines, we provide an approximate cost of obtaining a BubbleView click map per image using 10-15 participants. Table 7.7 contains a breakdown of costs that can be used as guidelines.

**Methodology limitations:** Compared to natural viewing or moving a mouse, clicking takes more time and effort, resulting in longer task timings and higher costs. Certain image regions which might not be as relevant to the task might never be clicked on, even though they may have received a quick glance in an eye tracking or moving-window setting. As a result, the image regions selected by clicks will tend to be more selective than the regions selected in these other settings. As shown in this chapter, the advantage of this selectivity is cleaner, more consistent results across participants. This can be used for determining the most important regions in an image (Exp. 4). But this comes at the potential disadvantage of certain image regions being missed, and other regions, like text, receiving disproportionate clicks (Exp. 1,3). How to encourage a more diverse sampling of image regions while maintaining all the other advantages of

**Table 7.7.** Total computed costs per image for obtaining the BubbleView clicks of 10–15 participants (both ends of the range included). These costs depend on how long, on average, participants spend on each image, which in turn depends on the task used. In the free-viewing setting, we fixed the time to either 10 or 30 seconds per image. In the description task, time is unconstrained, and participants move on to the next image after submitting their description for the previous image. During piloting, we estimated time per image for clicking and describing to take about 1.5 minutes. In reality, it took on average 3.2 minutes per image. The description task is more expensive but provides higher-quality click data and an additional data source: the descriptions themselves. These descriptions also serve as quality-control: the clicks of participants who generated poor-quality descriptions can be discarded.

| Task | Time/image | Images/HIT | Cost/HIT | Participants | Cost/image |
|------|-----------|-----------|---------|-------------|-----------|
| Free-viewing | 10 sec | 17 | $0.30 | 10–15 | $0.18–$0.26 |
| Free-viewing | 30 sec | 17 | $0.90 | 10–15 | $0.53–$0.79 |
| Description | 180 sec | 3 | $0.50 | 10–15 | $3.34–$5.00 |

BubbleView is a question for future investigations.

## ■ 7.7 Future directions

In this chapter we presented BubbleView, a mouse-contingent methodology to approximate eye fixations using mouse clicks. We validated BubbleView by conducting a series of experiments on different image stimuli and comparing clicks to eye fixations, importance maps, and mouse movements. We showed that BubbleView can reasonably approximate fixations, be used to collect image importance driven by human perception, and has a number of advantages compared to the moving-window approach, including better performance with fewer participants.

We analyzed BubbleView in the context of 4 image types (information visualizations, natural images, static webpages, and graphics designs), with 2 task types (free-viewing and description), with different task timings, image blur and bubble sizes, and different numbers of study participants. We provided the interested experimenter with some guidelines on how to use BubbleView for different tasks, how to select parameters, and which settings we found to work best under different conditions. Here we provide additional ideas of how BubbleView can be used and built on top of.

**Integrating BubbleView into crowdsourcing pipelines:** Unlike eye tracking experiments, BubbleView experiments can be feasibly ported online for the efficient and scalable collection of data using crowdsourcing. Large amounts of data call for data filtering and analysis methods that can scale as well. As shown in this chapter, BubbleView clicks can be analyzed automatically. In cases where text input is also collected from participants, filtering and analysis may require additional manual effort. However, it is possible to consider crowdsourcing pipelines where the data collected from the BubbleView tasks is piped directly into filtering tasks.

Following the idea of question-answering tasks, BubbleView can be incorporated into multi-player crowdsourcing games (e.g., ESP Game [216]). For instance, one participant

can generate questions, while the other participant answers using BubbleView clicks. In this setting, the first participant queries and supervises the responses of the second participant. In such a way both data collection and data cleaning can be built into the game.

**BubbleView data for training computational models:**   BubbleView can be used to generate large datasets for training computational models. BubbleView click maps on images can be used as importance maps for those images, and computational models can learn from this data to make predictions for new images. This could open up many interesting applications such as thumbnailing and retargeting designs or providing automatic design feedback (Chapter 10).

**Measuring information content:**   Clicking on an image region takes more effort than mousing over, and in turn, glancing at it. There is likely a relationship between the information content of an image region and the likelihood with which it is clicked, moused over, and glanced at. Clicking imposes a kind of energy barrier on the image content that will be explored by participants. Given a targeted task such as describing an image, participants are motivated to click in as few regions as necessary to reduce the overall effort and total task time. As a result, they tend to click in the most informative regions. Increasing the bubble size lowers this energy barrier: participants become less selective of where they're clicking when they can expose more of the image with each click. Changing the image blur also affects which image region will be clicked, based on its information content. More deeply studying the relationship between visual feature size, information content, image blur and bubble size is likely to provide some interesting insights. In the present study, by virtue of the images we selected for our experiments (e.g., to contain legible text) and the narrow range of image sizes we used, results were pretty stable across blur and bubble settings.

**Extending BubbleView to other tasks:** The interested experimenter may also choose to use BubbleView in settings and with parameters beyond the ones in this chapter, which leaves many possibilities for future investigation. For instance, Bubble-View can easily be extended to other visual attention tasks including visual search[6]. To implement a version of visual search using BubbleView, participants can be shown a blurred image and asked to find something in the image (e.g., an object in a natural scene, a specific piece of information in a graph, or an element in a graphic design). Task time can be either fixed, contingent on when the participant chooses to continue to the next image, or contingent on the participant's clicks (i.e., moving to the next image after the correct/expected location is clicked, or after a fixed number of clicks).

Another possible use for BubbleView is modifying the description task into a question-answering task. Participants can be asked to answer a specific question about the image by clicking around the blurred image to expose the content underneath. Each answer, correct, incorrect, or subjective, can be analyzed together with the sequence of clicks made (similar to Das et al. [45]).

---

[6]Some examples of visual attention tasks with operational definitions and recommended evaluations are included in Bylinskii et al. [32].

While we originally designed BubbleView as a more efficient alternative to collecting eye fixations on images, we have also shown in this chapter that it can be used to measure the importance of different image regions. This idea can be pushed even further in the future, using BubbleView to narrow in on image regions most useful for answering specific questions, extracting particular insights, or completing specific visual tasks. We showed that BubbleView generalizes to different types of images, including natural scenes, visualizations, websites, and graphic designs. This can be expanded to new image types, for instance for studying medical images, geographical maps, user interfaces, slides and posters. For future explorations, we provide our tool and code for launching experiments at massvis.mit.edu/bubbleview.

# ZoomMap: using zoom to capture user areas of interest on images

*This chapter is based on unpublished work by: Bylinskii, Z., Tancik, M., Newman, A., Zhong, K., Madan, S., Oliva, A., Durand, F. "ZoomMap: Using Zoom to Capture User Areas of Interest on Images"*

I N the last chapter we saw that by blurring images and having users click to deblur small bubble regions at full resolution, eye fixations can be well approximated for certain image types. *BubbleView*, the interface that was presented, is an implementation of the more general moving-window methodology, in which a limited amount of information is visible through a variable size window continuously following a cursor position [9, 101, 173] (see also Chapter 3.1).

In this chapter, we propose an extension to the moving-window methodology that takes advantage of the viewing behavior on a mobile screen. The mobile screen provides a naturally restricted window that is frequently used to explore multi-scale content with the help of the zoom functionality. We present an approach to capture zoom behavior using a mobile interface, a visualization of the spatial zoom patterns we call *ZoomMaps* that can be overlaid on images, and applications that can be built with zoom data. As an under-explored modality, we show that zoom can be used to measure the natural scale of visual content, approximate image saliency, act as a debugging tool for designs, and be harnessed for applications like personalized thumbnailing.



**Figure 8.1.** We built an image gallery website with tracking capabilities to capture how people use the zoom functionality during image viewing.

**Figure 8.2.** (a) A poster that is zoomed in using the mobile interface at the region outlined in orange. (b) The zoom level for an image region is computed as the quotient of the full image area divided by the area of the image region that has been magnified. (c) A ZoomMap for an image visualizes the average zoom of each pixel over an entire viewing period. (d) We show the zoom profiles of 3 pixels sampled from this zoom map: plotted are the zoom levels of each pixel over time. (e) The final ZoomMap overlaid on top of the original poster image.

## ■ 8.1  User interface for capturing zoom

To capture mobile zoom information, we built an image gallery website with tracking capabilities. The goal was to use an interface familiar to mobile device users. The gallery was based on the PhotoSwipe Javascript library[1], modified to capture and store any changes to the visible region of the image along with a timestamp. The website allows pinching to zoom into an image, and swiping to switch images (Fig. 8.1). When the image is swiped, all the interaction data on the image is stored to a server. The interaction data contains the coordinates of the zoomed-in portion of the image with a timestamp for every *event* triggered by the user (i.e., the CSS rescales or repositions the image on the screen).

## ■ 8.2  ZoomMaps: visualizing spatial zoom patterns

From our mobile interface, we have the X and Y coordinates of zoomed image regions within the coordinate frame of the whole image, along with timestamps of when regions

---

[1]https://github.com/dimsemenov/photoswipe

were in focus. We can use this information to extract (1) image regions viewed the longest, and (2) image regions viewed under different zoom levels.

To display spatial zoom patterns on images, we construct the following visualization: for every pixel in the image, we compute its average *zoom level* over the entire viewing interval. We define the zoom level for an image region as the quotient of the full image area divided by the area of the image region that has been magnified (Fig. 8.2a). We assign this zoom level to all the pixels contained in the image region. By aggregating these zoom levels over time, we compute each pixel's average zoom level. We can then visualize this value per pixel to obtain our *ZoomMaps* (Fig. 8.2b). Higher values in the ZoomMaps correspond to regions of the image that were inspected with closer zoom on average.

## ■ 8.3 Experiments with natural images

### ■ 8.3.1 Procedure

**Stimuli:** We used images from the CAT2000 database [20], a dataset with ground truth eye tracking data used for evaluating saliency models on the MIT Saliency Benchmark [31]. CAT2000 contains images from 20 image categories. We sampled 5 images each from: *Action, Art, Fractal, Outdoor Man Made, Outdoor Natural, Satellite, Social*. These 7 categories were chosen to capture a diversity of image content and to be viewable at different scales.

**Participants:** Recruitment was done through a departmental mailing list. We obtained 14 participants (6 females), including computer science students and staff (ages 19-29). Participants were paid $5 each.

**Task:** The task consisted of viewing a gallery of 35 images in landscape orientation on a mobile screen. Images were randomly shuffled for each participant. Participants were asked to spend 10 minutes to *"explore each image carefully, using the pinch zoom gestures of the phone"*. They were also told that upon completion, they would be sent a questionnaire to test their memory for image details. The main purpose of these instructions was to encourage careful viewing and attention to the task. Participants were then sent a brief post-viewing questionnaire to collect their demographic info, feedback about the task, and ask them to (1) describe 5 images from memory, and (2) describe how they used the zoom functionality to view the images. Answers to (1) are not analyzed in this work. Answers to (2) are discussed below.

### ■ 8.3.2 Results

**What do people zoom into?** Participants self-reported on regions of images they zoomed into during the study. We tallied the number of times different image elements were mentioned in the responses (Fig. 8.3). Image elements that people zoom are also reported to be frequently fixated in eye tracking studies [21, 33, 38]. As discussed in Chapter 6.2 the regions attracting the most eye fixations (across the MIT300 dataset of natural images) include *people*, *text*, and *faces*, and often contain *unusual objects*, the

What do people zoom into?

**Figure 8.3.**  (Right) Self-reporting of participants about where they choose to zoom in images. This histogram plots the number of mentions of each of the elements from 14 participant questionnaires.  See also Table 6.2 from Chapter 6.

*main subject* of the image, and *animals*. In the present study, we also had fractal, art, and satellite images that fall into different categories than the natural images previously analyzed. A third of our participants explicitly mentioned zooming the aerial images. Most participants reported that they did not zoom fractal images.

To verify these claims, we computed the average zoom level over all images per category. We found that people zoomed most in satellite (1.88) and art (1.85) images, and least in fractal (1.29) and outdoor natural (1.32) images (average zoom levels in parentheses). About 25% of participants did not even use the zoom functionality on the fractal and outdoor natural images. We note that faces and people were distributed across most of the image categories.

**How well do ZoomMaps correlate with image saliency?** Zoom provides a coarse notion of attention, since an observer's gaze can move around the mobile screen without additional interactions, and only when a particular image region is both (i) interesting to the observer and (ii) not fully visible at the current screen resolution, will the zoom functionality be invoked. As a result, we can expect at most a couple of zooming actions on the average photograph (contrary to more dense content such as the posters in the next section). Individual participant ZoomMaps on the CAT2000 images were quite blocky, with 1-2 regions of zoom (if any). However, by averaging the ZoomMaps across multiple participants (8-14 in our case), we obtain smoother maps with clearer areas of interest (AOIs). These AOIs correspond to image regions that are salient, as measured by eye movements (Fig. 8.4). The CAT2000 images are available with ground truth fixation maps, obtained by aggregating and smoothing the eye fixation locations of 24 observers.

We compared our average ZoomMaps ($Z$) to the fixation maps ($F$) on 35 CAT2000 images by computing Pearson's Correlation Coefficient, CC (see Chapter 5, Sec. 5.4.2 for the computation). We obtained an average CC score of 0.51 across the 35 images. To put this into perspective, a center prior map (a central Gaussian) achieves a CC score of 0.46 when compared to the fixation maps, while the IO (inter-observer) model achieves

|  | Images Viewed | Individual ZoomMaps | | | Avg. ZoomMaps | Fixation Maps |



**Figure 8.4.** Participants viewing images on a mobile screen zoom into a few areas of interest, leading to coarse maps of attention. Averaging ZoomMaps across 14 participants yields smoother maps that can be used to approximate image saliency (eye fixation maps). In both cases, people either look or zoom into regions they want to explore further.

0.56 (chance is at 0). The IO model is computed by comparing the fixation map of $N-1$ observers to the fixation map of the remaining observer, and then averaging over all observers. So while ZoomMaps can not fully account for saliency, they are highly correlated with it, and can serve as a coarse approximation for some applications, including image compression and transmission (see *Discussion*).

## ■ 8.4 Experiments with academic posters

### ■ 8.4.1 Procedure

**Stimuli:** We obtained 13 academic posters from our colleagues from the 2017 conference on *Computer Vision and Pattern Recognition*. Because they were prepared for the same conference, the posters were all in landscape orientation and had similar dimensions. We resized them to approximately $3000 \times 1500$ pixels, while maintaining the original aspect ratios.

**Participants:** We recruited 11 undergraduate and graduate students (3 females) who had familiarity with computer vision, and paid them $25 each. Three participants were excluded after the completion of the first part of the experiment for not following instructions completely.

**Task:** Participants completed a task split across 2 days. In the first part, they viewed 6 posters for 30 minutes total. They were told to carefully study each poster because their memory for the posters would be tested. On the next day, participants were asked the same 4 questions about all 6 posters (about contributions, results, applications, and key concepts). After answering all the questions for a poster, they were asked to rate *"how well do you remember this poster?"* on a scale from 1 (*not at all*) to 5 (*quite clearly*). Since we believe these numerical values reflect how participants answered the previous questions, we use these values for evaluation (verbal answers

are harder to quantify). Key to this part of the study is that prior to answering the questions for each poster, participants were shown a thumbnail of the poster, roughly 10% of its original size, for 10 seconds. This size was too small to be able to read any text, but was intended to remind the participants of a particular poster. For a total of 6 posters, they saw 2 static thumbnails, 2 GIFs computed from their own zoom patterns, and 2 other GIFs (1 computed from another user's zoom, and 1 random baseline). At the end of the study participants were asked to indicate whether they preferred the static thumbnails or GIFs for recalling poster content.

### ■ 8.4.2  Computing GIFs

Using the zoom data captured by our mobile interface, we obtain the image regions viewed the longest by each participant. These windows could come from any part of the poster, and may be at different zoom levels (Fig. 8.6, top). However, they all have the same aspect ratio since they filled the mobile screen size. We take the 5 longest-viewed image regions and concatenate them into a GIF of 6 frames (800 ms each), where the first frame is of the whole poster. In 10 seconds, the GIF played twice. Note that an alternative future design could zoom-out, pan, and zoom-in between frames to provide more context.

### ■ 8.4.3  Results

**Using zoom data for thumbnailing:** During viewing, participants zoom a region of an image to bring it into focus and study it more carefully. By snapshotting these image regions, we can obtain a temporal sequence of "frames" that we can play back. We used this idea to create small GIFs (not unlike YouTube video previews) of each user's viewing patterns for each poster. During the survey at the end of the study, 6/8 participants indicated they preferred seeing these GIFs rather than the static thumbnails for recalling poster content because: the GIFs *"zoom into the image, enabling the viewer to see words and figures up close", "it was easier to see details", "it was nice to see the images closer"* (participant comments). Two of the participants remarked that it helped to see the views of the poster they had originally looked at.

We also considered participant ratings of how well each poster was remembered in 3 conditions: when shown static thumbnails, GIFs of their own zoom patterns, and GIFs of other zoom patterns. Participants were shown 2 of each type, so we averaged their ratings across both instances. Recall that participants were shown each thumbnail at the *beginning* of the questionnaire, then had to rate how well they remembered each poster at the *end* of the questionnaire (only after answering questions about the poster). Participants better remembered posters for which they were shown one of the GIF thumbnails over the static thumbnails in 7/8 of the cases (Fig. 8.5). In the remaining case, posters were remembered equally well for both conditions. We did not find evidence that participants better remembered posters for their own versus other users GIFs, but this aspect requires further investigation.

**Rating of memory after viewing thumbnail**

**Figure 8.5.** (Left) How well participants self-reported they remembered the content of posters after seeing a thumbnail that was either static, or a GIF - of their own or someone else's viewing patterns, bringing different image content to the forefront of the thumbnail.

**How consistent are viewing patterns on posters?** We measured the consistency between participant zoom patterns by computing the inter-observer model for all the poster images (comparing the ZoomMap of $N-1$ participants to the ZoomMap of the remaining participant). We obtained a lower CC score of 0.41 (compared to 0.63 among the CAT2000 images). Qualitatively examining the ZoomMaps of individuals, there is a lot more variability than with photographs (Fig. 8.6). Where people choose to focus in academic posters depends on their personal background knowledge and interests. General models of saliency for posters may be less appropriate, but personalized applications (e.g., personalized thumbnails) may be promising.

**What can zoom tell us about poster design?** Academic posters contain a lot of multi-scale information, and need to be zoomed to be viewed on a mobile device. The average zoom level across all poster images[2] and participants was 4.15 (compared to 1.57 for the CAT2000 images). This means that image regions were viewed at over 4 times the resolution they originally appear at on the mobile screen. The largest average zoom level for an individual poster was 5.11 and the smallest was 3.23 (Fig. 8.7). Given that all posters were presented at the conference at the same size, the mobile viewing data hints at the fact that the poster with the largest zoom level had a scaling issue: information was presented too small for comfortable viewing. Our mobile interface can be viewed as a debugging tool for designs, providing an evaluation of the appropriateness of the physical scale at which information is visually presented.

## ■ 8.5  Future directions

The use of the zoom functionality during image viewing on mobile devices has been under-explored as a tool to study user viewing behavior, attention, and interest. We presented an initial investigation into how zoom can be recorded, visualized, and utilized to make inferences about visual content (e.g., predict saliency on photographs and debug the design of posters). Here we highlight the potential applications that can be built

---

[2]For a more robust evaluation, we only include the 8 posters with at least 7 viewers each.

Individual Zoom Patterns



Individual ZoomMaps



**Figure 8.6. Top:** the 5 longest-viewed image regions by an individual participant. Note that they occur at different zoom levels. We would sequence these 5 regions into a personalized GIF. **Bottom:** a sample of 3 individual ZoomMaps for the same poster, showing diversity in the regions people choose to explore.



**Figure 8.7.** (a) Poster with higher than average zoom level. (b) Poster with lower than average zoom level. In the top poster, images and text are smaller (compared to the bottom poster) and participants had to zoom more to see the content better. As a debugging tool, viewing posters on a mobile screen can help identify how people would view the poster when it is printed.

on top of zoom data.

**Image compression and transmission:** We showed that participants zoom more on specific content in photographs (e.g., people, faces, text) and on certain types of images (e.g., aerial) over others (e.g., natural outdoor scenes). This information can be encoded into image applications, such as photo galleries, to preload in high resolution only what will be predictably zoomed. Conversely, images can be compressed in a way where detail is preserved in regions of the image that are more likely to be attended to. Image compression and transmission have long been touted as potential applications of saliency. Zoom information can provide another, in some ways more direct, measurement of the resolution requirements for different image regions.

**Personalized thumbnails:** Because mobile devices are so ubiquitous and used on a daily basis, there are many opportunities for personalization. We demonstrated that zoom data can be converted into personalized GIFs of viewing patterns, which may help to record and effectively remind a user of information that was previously studied. Extensions include snapshotting and thumbnailing articles, e-mails, and images in the spots where users studied them most closely, as place-holders or memory hooks.

**Automatic design feedback:** The amount different observers zoom into different portions of a graphic design (e.g., poster, website) can provide important insights about the quality of the design. If all observers zoom into the design more than other designs (of the same target physical size), perhaps the design content is inappropriately scaled for the target application. If all observers are zooming into the same region of the design, perhaps that region needs to be resized relative to the rest of the design. Since zoom can be used to infer the natural scale of different regions of a design, this can be built into an automatic redesign tool that aims to rebalance content to an appropriate scale.

# Part IV

# Predicting Attention on Visualizations and Graphic Designs

# Chapter 9

# Learning visual importance

*This chapter is based on: Bylinskii, Z., Kim, N.W., O'Donovan, P., Alsheikh, S., Madan, S., Pfister, H., Durand, F., Russell, B., Hertzmann, A. "Learning Visual Importance for Graphic Designs and Data Visualizations", ACM User Interface Software and Technology Symposium (UIST 2017).*

**A**CRUCIAL goal of any graphic design or information visualization is to communicate the relative *importance* of different design elements, so that the viewer knows where to focus attention and how to interpret the design. In other words, the design should provide an effective management of attention [180]. Understanding how viewers perceive a design could be useful for many stages of the design process; for instance, to provide feedback [185]. Automatic understanding can help build tools to search, retarget, and summarize information in designs and visualizations. Though saliency prediction in natural images has recently become quite effective, there is little work in importance prediction for either graphic designs or information visualizations.

We use **importance** as a generic term to describe the perceived relative weighting of design elements. Image saliency, which has been studied extensively, is a form of



**Figure 9.1.** We present two neural network models trained on crowdsourced importance. We trained the graphic design model using a dataset of 1K graphic designs with GDI annotations [154]. For training the information visualization model, we collected mouse clicks using the BubbleView methodology [113] on 1.4K MASSVIS information visualizations [23]. Both networks successfully predict ground truth importance and can be used for applications such as retargeting, thumbnailing, and interactive design tools. Warmer colors in our heatmaps indicate higher importance.

**Figure 9.2.** We show an interactive graphic design application using our model that lets users change and visualize the importance values of elements. Users can move and resize elements, as well as change color, font, and opacity, and see the updated realtime importance predictions. For instance, a user changes the color of the text to the left of the runner to increase its importance (middle panel). The rightmost panel includes a few additional changes to the size, font, and placement of the text elements to modify their relative importance scores. A demo is available at `visimportance.csail.mit.edu`.

importance. However, whereas traditional notions of saliency refer to bottom-up, pop-out effects, our notion of importance can also depend on higher-level factors such as the semantic categories of design elements (e.g., title text, axis text, data points).

Recently, deep learning methods, trained on large datasets, have produced a substantial jump in performance on standard saliency benchmarks (Chapter 6). However, these methods have been developed exclusively for analyzing natural images, and are not trained or tested on graphic designs. The work presented in this chapter is the first to apply neural network importance predictors to both graphic designs and information visualizations. We use a state-of-the-art deep learning architecture, and train models on two types of crowdsourced importance data: graphic design importance (GDI) annotations [154] and a dataset of clicks we collected on information visualizations using the BubbleView interface (Chapter 7).

Our importance models take input designs in bitmap form. The original vector data is not required. As a result, the models are agnostic to the encoding format of the image and can be applied to existing libraries of bitmap designs. Our models pick up on some of the higher-level trends in ground truth human annotations. For instance, across a diverse collection of visualizations and designs, our models learn to localize the titles and correctly weight the relative importance of different design elements (Fig. 9.1).

In Chapter 10 we show how the predicted importance maps can be used as a common building block for a number of different applications, including retargeting and thumbnailing. Our predictions become inputs to cropping and seam carving with almost no additional post-processing. Despite the simplicity of the approach, our retargeting and thumbnailing results are on par with, or outperform, related methods, as validated by a set of user studies launched on Amazon's Mechanical Turk (MTurk). Moreover, an advantage of the fast test-time performance of neural networks makes it feasible for our predictions to be integrated into interactive design tools (Fig. 9.2). With another set of

**Figure 9.3.** Left: Comparison of eye movements collected in a controlled lab setting [24], and clicks that we crowdsourced using the BubbleView interface [112, 113]. Right: Comparison of importance annotations from the GDI dataset, and clicks that we crowdsourced using the BubbleView interface. These examples were chosen to demonstrate some of the similarities and differences between the modalities. For instance, compared to eye fixations, clicks are sometimes more concentrated around text. Compared to the GDI annotations, clicks do not assign uniform importance to whole design elements. Despite these differences, BubbleView data leads to similar importance rankings of visual elements (*Evaluation*).

user studies, we validate that our model generalizes to fine-grained design variations and correctly predicts how importance is affected by changes in element size and location on a design.

**Contributions:** We present two neural network models for predicting importance: in graphic designs and information visualizations. This is the first time importance prediction is introduced for information visualizations. For this purpose, we collected a dataset of BubbleView clicks on 1,411 information visualizations. We also show that BubbleView clicks are related to explicit importance annotations [154] on graphic designs. We collected importance annotations for 264 graphic designs with fine-grained variations in the spatial arrangement and sizes of design elements. We demonstrate how our importance predictions can be used for retargeting and thumbnailing, and include user studies to validate result quality. Finally, we provide a working interactive demo.

## ■ 9.1 Data Collection

To train our models we collected BubbleView data [112, 113] for information visualizations, and used the Graphic Design Importance (GDI) dataset by O'Donovan et al. [154] for graphic designs. We compared different measurements of importance: BubbleView clicks to eye movements on information visualizations, and BubbleView clicks to GDI annotations on graphic designs.

### ■ 9.1.1  Ground truth importance for information visualizations

Large datasets are one of the prerequisites to train neural network models. Unfortunately, collecting human eye movements for even hundreds of images is extremely expensive and time-consuming. Instead, we use the BubbleView interface introduced in Chapter 7 to collect clicks on images as a proxy to eye fixations. Unlike eye tracking, which requires expensive equipment and a controlled lab study, BubbleView can be used to to collect large datasets with online crowdsourcing.

Concurrent work in the computer vision community has applied a similar methodology to natural images. SALICON [103] is a crowdsourced dataset of mouse movements on natural images that has been shown to approximate free-viewing eye fixations. Current state-of-the-art models on saliency benchmarks have all been trained on the SALICON data [41, 93, 129, 157, 231]. BubbleView was concurrently developed [112] to approximate eye fixations on information visualizations with a description task. Some advantages of BubbleView over SALICON are discussed in Chapter 7 (Sec. 7.4.2).

Using Amazon's Mechanical Turk (MTurk), we collected BubbleView data on a set of 1,411 information visualizations from the MASSVIS dataset [23], spanning a diverse collection of sources (news media, government publications, etc.) and encoding types (bar graphs, treemaps, node-link diagrams, etc.). We manually filtered out visualizations containing illegible and non-English text, as well as scientific and technical visualizations containing too little context. Images were scaled to have a maximum dimension of 600 pixels to a side while maintaining their aspect-ratios to fit inside the MTurk task window. We blurred the visualizations using a Gaussian filter with a radius of 40 pixels and used a bubble size with a radius of 32 pixels as in [113]. MTurk participants were additionally required to provide descriptions for the visualizations to ensure that they meaningfully explored each image. Each visualization was shown to an average of 15 participants. We aggregated the clicks of all participants on each visualization and blurred the click locations with a Gaussian filter with a radius of 32 pixels, to match the format of the eye movement data.

We used the MASSVIS eye movement data for testing our importance predictions. Fixation maps were created by aggregating eye fixation locations of an average of 16 participants viewing each visualization for 10 seconds. Fixation locations were Gaussian filtered with a blur radius of 32 pixels. Fig. 9.3a includes a comparison of the BubbleView click maps to eye fixation maps from the MASSVIS dataset.

### ■ 9.1.2  Ground truth importance for graphic designs

We used the Graphic Design Importance (GDI) dataset [154] which comes with importance annotations for 1,078 graphic designs from Flickr. Thirty-five MTurk participants were asked to label important regions in a design using binary masks, and their annotations were averaged. Participants were not given any instruction as to the meaning of "importance." To determine how BubbleView clicks relate to explicit importance annotations, we ran the BubbleView study on these graphic designs and collected data

from an average of 15 participants per design. Fig. 9.3b shows comparisons between the GDI annotations and BubbleView click maps. In both data similar elements and regions of designs emerge as important.

Each representation has potential advantages. The GDI annotations assign a more uniform importance score to whole elements. This can serve as a soft segmentation to facilitate design applications like retargeting. BubbleView maps may be more appropriate for directly modeling human attention.

## ■ 9.2 Models for predicting importance

Given a graphic design or information visualization, our task is to predict the importance of the content at each pixel location. We assume the input design/visualization is a bitmap image. The output importance prediction at each pixel $i$ is $P_i \in [0,1]$, where larger values indicate higher importance. We approach this problem using deep learning, which has led to many recent breakthroughs on a variety of image processing tasks in the computer vision community [124, 175], including the closely related task of saliency modeling.

Similar to some top-performing saliency models for natural images [93, 125], our architecture is based on fully convolutional networks (FCNs) [141]. FCNs are specified by a directed acyclic graph of linear (e.g., convolution) and nonlinear (e.g., max pool, ReLU) operations over the pixel grid, and a set of parameters for the operations. The network parameters are optimized over a loss function given a labeled training dataset. We refer the reader to Long et al. [141] for more details.

We predict real-valued importance using a different training loss function from the original FCN work, which predicted discrete object classes. Given ground truth importances at each pixel $i$, $Q_i \in [0,1]$, we optimize the sigmoid cross entropy loss for FCN model parameters $\Theta$ over all pixels $i = 1, \ldots, N$:

$$L(\Theta) = -\frac{1}{N} \sum_{i=1}^{N} \left( Q_i \log P_i + (1 - Q_i) \log(1 - P_i) \right) \tag{9.1}$$

where $P_i = \sigma\left(f_i(\Theta)\right)$ is the output prediction of the FCN $f_i(\Theta)$ composed with the sigmoid function $\sigma(x) = (1 + \exp(-x))^{-1}$. Note that the same loss is used for binary classification, where $Q_i \in \{0, 1\}$. Here, we extend it to real-valued $Q_i \in [0, 1]$. We use a different loss than other saliency models based on neural networks that optimize Euclidean [125, 157], weighted Euclidean [41], or binary classification losses [129, 231]. Our loss is better suited to $[0, 1]$ values, and is equivalent to optimizing the KL loss commonly used for saliency evaluation.

We trained separate networks for information visualizations and for graphic designs. For the information visualizations, we split the 1.4K MASSVIS images for which we collected BubbleView click data into 1,209 training images and 202 test images. For the test set we chose MASSVIS images for which eye movements are available [24]. For the graphic designs, we split the 1,078 GDI images into 862 training images and 216 test

**Figure 9.4.** We increase the precision of our FCN-32s predictions by combining output from the final layer of the network with outputs from lower levels. The resulting predictions, FCN-16s and FCN-8s, capture finer details. We found FCN-16s sufficient for our model for graphic designs, as FCN-8s did not add a performance boost. For our model for information visualizations, we found no performance gains beyond FCN-32s.

images (80-20% split). We used the GDI annotations [154] for training. We found that training on the GDI annotations rather than the BubbleView clicks on graphic designs facilitated the design applications better, since the GDI annotations were better aligned to element boundaries.

**Model details:** We converted an Oxford VGG-16 convolutional neural network [198] to an FCN-32s model via network surgery using the implementation in Caffe [102]. The model's predictions are 1/32 of the input image resolution, due to successive pooling layers. To increase the resolution of the predictions and capture fine details, we followed the procedure in Long et al. [141] to add skip connections from earlier layers to form FCN-16s and FCN-8s models, that are respectively, 1/16 and 1/8 of the input image resolution. We found that the FCN-16s (with a single skip connection from *pool4*) improved the graphic design importance maps relative to the FCN-32s model (Fig. 9.4), but that adding an additional skip connection from *pool3* (FCN-8s) performed similarly. We found that skip connections lead to no gains for the information visualization importance. For our experiments we used the trained FCN-16s for graphic designs and the FCN-32s for information visualizations.

Since we have limited training data we initialized the network parameters with the pre-trained FCN32s model for semantic segmentation in natural images [141], and fine-tuned it for our task. The convolutional layers at the end of the network and the skip connections were randomly initialized.

We opted for a smaller architecture with fewer parameters than some other neural network saliency models for natural images. This makes our model more effective for our datasets, which are currently an order-of-magnitude smaller than the natural image saliency datasets.

**Training details:** The FCN-32s network was initialized with a base learning rate ($lr$) of $1e-05$, scaled by a factor of 0.1 every 20K iterations. A stochastic gradient descent [25] solver with a momentum of 0.9 and weight decay of 0.0005 was used, and

| Model | CC score ↑ | KL score ↓ |
|---|---|---|
| Chance | 0.00 | 0.75 |
| Judd [106] | 0.11 | 0.49 |
| DeepGaze [129] | 0.57 | 3.48 |
| **Our model** | **0.69** | **0.33** |

**Table 9.1.** How well can our importance model predict the BubbleView click maps? We add comparisons to two other top-performing saliency models and a chance baseline. Scores are averaged over 202 test information visualizations. A higher CC score and lower KL score are better.

run for 100K iterations. The FCN-16s network was initialized with the weights of the FCN-32s network and a base $lr$ of $1e-11$ (equivalent to the learning rate used on the last iterations of training the FCN-32s network, and scaled by 0.001). The rest of the training parameters were the same. The FCN-8s network was similarly initialized with the weights of the FCN-16s network and a base $lr$ of $1e-17$. Our learning rate schedule was similar to the one used for semantic segmentation [141].

## ■ 9.3  Evaluation of model predictions

We compare the performance of our two importance models to ground truth importance on each dataset. For information visualizations, we compare predicted importance maps to bubble clicks gathered using BubbleView, and to eye fixations from the MASSVIS dataset. For graphic designs, we compare predicted importance maps to GDI annotations. For evaluation, we use the Kullback-Leibler divergence (KL) and Pearson's Correlation Coefficient (CC) metrics (see Chapter 5). For further intuition about how KL and CC metrics score similarity, we provide scores above each image in Fig. 9.5, showing high- and low-scoring predictions.

### ■ 9.3.1  Prediction performance on information visualizations

We include predictions from our importance model in Fig. 9.5. Notice how we correctly predict important regions in the ground truth corresponding to titles, captions, and legends. We quantitatively evaluate our approach on our collected dataset of Bubble-View clicks. We report CC and KL scores averaged over our dataset of 202 test images in Table 9.1.

   We compare against the following baselines: chance, Judd saliency [106], and Deep-Gaze [129], a top neural network saliency model trained on the SALICON dataset [103] of mouse movements on natural images. The chance baseline, used in saliency benchmarks [36, 107], is computed by uniformly sampling a real value between 0 and 1 at each image pixel. Our approach out-performs all baselines. KL is highly sensitive to false negatives and drastically penalizes sparser models (Chapter 5.5), explaining the high KL values for DeepGaze in Table 9.1. Post-processing or directly optimizing models for specific metrics can yield more favorable performances [131]. Because of the sensitivity of KL to output regularization, we advise against using it (solely) to compare models.

   How well does our neural network model, trained on clicks, predict eye fixations?

| Model | CC score ↑ | KL score ↓ |
|---|---|---|
| Chance | 0.00 | 1.08 |
| Judd [106] | 0.19 | 0.74 |
| DeepGaze [129] | 0.53 | 3.10 |
| **Our model** | **0.54** | **0.63** |
| Bubble clicks | 0.79 | 0.28 |

**Table 9.2.** How well can human eye fixations be predicted? We measured the similarity between human fixation maps and various predictors. Scores are averaged over 202 test information visualizations. A good model achieves a high CC score and low KL score. Our neural network model was trained on BubbleView click data, so that is the modality it can predict best. Nevertheless, its predictions are also representative of eye fixation data. As an upper bound on this prediction performance, we consider how well the BubbleView click data predicts eye fixations, and as a lower bound, how well chance predicts eye fixations.

We find that the predicted importance is representative of eye fixation patterns as well (Table 9.2), although the difference in scores indicates that our model might be learning from patterns in the click data that are different from fixations.

**Which elements are most important?** For our analysis, we used the element segmentations available for the visualizations in the MASSVIS dataset [24]. We overlapped these segmentations with normalized maps of eye fixations, clicks, and predicted importance. We computed the max score of the map within each element to get an importance ranking across elements (we extended the analysis in Chapter 7, Sec. 7.3.1 to include the predicted importance maps; see also Fig. 7.7).

Text elements, such as titles and captions, were the most looked at, and clicked on, elements, and were also predicted most important by our model (Fig. 9.6). Even though our model was trained on BubbleView clicks, the predicted importance remains representative of eye fixation patterns. With regards to differences, our model overpredicts the importance of titles. Our model learns to localize visualization titles very well (Fig. 9.5). This matches findings on human perception (Chapter 4), as among the text and other content in a visualization, titles tend to be best remembered by human observers [24].

## ■ 9.3.2 Prediction performance on graphic designs

The closest approach to ours is the work of O'Donovan et al. [154] who computed an importance model for the GDI dataset. We re-ran their baseline models on the train-test split used for our model (Table 9.3). To replicate their evaluation, we report root-mean-square error (RMSE) and the $R^2$ coefficient, where $R^2 = 1$ indicates a perfect predictor, and $R^2 = 0$ is the baseline of predicting the mean importance value. Defining $Q$ as the ground truth importance map and $P$ as the predicted importance map, we iterate over all pixels $i$ to compute:

**Figure 9.5.** Importance predictions for information visualizations, compared to ground truth BubbleView clicks and sorted by performance. Our model is good at localizing visualization titles (the element clicked on, and gazed at, most by human participants) as well as picking up the extreme points on graphs (e.g., top and bottom entries). We include a failure case where our model overestimates the importance of the visual map regions.



**Figure 9.6.** Relative importance scores of different elements in an information visualization assigned by eye fixation maps, BubbleView click maps, and model predictions. Scores were computed by overlapping element segmentations with normalized importance maps, and taking the max of the map within each element as its score. The elements that received the most clicks also tended to be highly fixated during viewing (Spearman's $r_s = .96, p < .001$). Text (titles, labels, paragraphs) received a lot of attention. Our neural network model correctly predicted the relative importance of these regions relative to eye movements ($r_s = .96, p < .001$).

$$RMSE(P,Q) = \sqrt{\frac{1}{N}\sum_{i=1}^{N}(Q_i - P_i)^2} \tag{9.2}$$

$$R^2(P,Q) = 1 - \frac{\sum_i (Q_i - P_i)^2}{\sum_i (Q_i - \overline{Q})^2} \quad \text{where } \overline{Q} = \frac{1}{N}\sum_{i=1}^{N} Q_i \tag{9.3}$$

The full O'Donovan model (*OD-Full*) requires manual annotations of *text*, *face*, and *person* regions, and would not be practical in an automatic setting. For a fair comparison, we evaluate our automatic predicted importance model (*Ours*) against the automatic portion of the O'Donovan model, which does not rely on human annotations (*OD-Automatic*). We find that our model outperforms *OD-Automatic*. Our model is also 100X faster, since it requires a single feed-forward pass through the network ($\sim$0.1 s/image on a GPU). O'Donovan's method requires separate computations of multiple CPU-based saliency models and image features ($\sim$10 s/image at the most efficient setting).

In Table 9.3, we include the performance of *Ours+OD*, where we added our importance predictions as an additional feature during training of the O'Donovan model, and re-estimated the optimal weights for combining all the features. *Ours+OD* improves upon *OD-Full* indicating that our importance predictions are not fully explainable by the existing features (e.g., text or natural image saliency). This full model is included for demonstration purposes only, and is not practical for interactive applications.

We also annotated elements in each of the test graphic designs using bounding boxes, and computed the maximum importance value in each bounding box as the element's score (Fig. 9.7). This is the same as the analysis in Chapter 7 (Sec. 7.4.1; see also Fig. 7.12), but extended to predicted importance maps. We obtain an average Spearman rank correlation of 0.56 between the predicted and ground truth scores assigned to the graphic design elements.

Some examples of predictions are included in Fig. 9.8. Our predictions capture important general trends, such as larger and more central text and visual elements being more important. However, text regions are not always well segmented (predicted importance is not uniform over a text element), and text written in unusual fonts is not always detected. Such problems could be ameliorated through training on larger datasets. Harder cases are directly comparing the importance of a visual and text, which can depend on the semantics of the text itself (how informative it is) and the quality of the visual (how unexpected, aesthetic, etc.).

### ■ 9.3.3 Prediction performance on fine-grained design variations

To check for feasibility of an interactive application we perform a more fine-grained test. We want the importance rankings of elements to be adjusted accurately when the user makes changes to their current design. For example, if the user makes a text box larger, then its importance should not go down in the ranking. Our predicted

| Model | RMSE $\downarrow$ | $R^2 \uparrow$ |
|---|---|---|
| Saliency | .229 | .462 |
| OD-Automatic | .212 | .539 |
| **Ours** | **.203** | **.576** |
| Annotations | .195 | .608 |
| Ours+Annot | .164 | .725 |
| OD-Full | .155 | .754 |
| **Ours+OD** | **.150** | **.769** |

**Table 9.3.** A comparison of our predicted importance model (*Ours*) with the model of O'Donovan et al. [154]. Lower *RMSE* and higher $R^2$ are better. Our model outperforms the fully automatic O'Donovan variant (*OD-Automatic*). Another fully automatic model from [154] is *Saliency*, a learned combination of 4 saliency models: Itti&Koch [97], Hou&Zhang [86], Judd et al. [106], and Goferman et al. [64]. We also report the results of the semi-automatic OD-Full model, which includes manual annotations of *text*, *face*, and *person* regions. The performance of these features is also reported separately as *Annotations*. When we combine our approach with OD-Full (*Ours+OD*), we can improve upon the OD model. Note that the first 3 rows of this table correspond to fully automatic models, while the last 4 include manual annotations. The top-performing model is bolded in each case.



**Figure 9.7.** An example comparison between the predicted importance of design elements and the ground truth GDI annotations. The heatmaps are overlapped with element bounding boxes and the maximum score per box is used as the element's importance score (between 0 and 1).

**Figure 9.8.** Importance predictions for graphic designs, sorted by performance. Performance is measured as the Spearman rank correlation (R) between the importance scores assigned to design elements by ground truth (GDI annotations) and predicted importance maps. A score of 1 indicates a perfect rank correlation; a negative score indicates the element rankings are reversed. The predicted importance maps distribute importance between text and visual features. We include a failure case where the importance of the man in the design is underestimated.

importance model has not been explicitly trained on systematic design variations, so we test if it can generalize to such a setting.

We used the Design Improvement Results dataset [154] containing 11 designs with an average of 35 variants. Across variants, the elements are preserved but the location and scale of the elements varies. We repeated the MTurk importance labeling task of O'Donovan et al. [154] on a subset of 264 design variants, recruiting an average of 19 participants to annotate the most important regions on each design. We averaged all participant annotations per design to obtain ground truth importance heatmaps. We segmented each design into elements and used the ground truth and predicted importance heatmaps to assign importance scores to all the elements, calculating the maximum heatmap value falling within each segment. The predicted and ground truth importance scores assigned to these elements achieved an average Spearman's correlation $r_s = .53$. As Fig. 9.9 shows, even though we make some absolute errors, we successfully account for the impact of design changes such as the location and size of various elements. Crucially, our model was not trained on systematic design variations, like changes in font, text size, or element location; nevertheless, it can correctly assign relative importance values to different design elements, as they are moved around and resized. This provides evidence that our model can provide meaningful predictions within an interactive tool setting (Chapter 10.4).

**Figure 9.9.** Sample input designs, and how the relative importance of the different design elements changes as they are moved, resized, and otherwise modified. For instance, compared to in (a), the event date stands out more and gains importance when it occurs at the bottom of the poster, in large font, on a contrasting background (b). Similarly, when the most important text of the design in (c) is moved to the upper righthand corner where it is not surrounded by other text, it gains prominence (d). Our automatic model makes similar predictions of the relative importance of design elements as ground truth human annotations.

### ■ 9.3.4  Limitations

Our neural network model is only as good as the training data we provide it. In the case of information visualizations, there is a strong bias, both by the model and the ground truth human data, to focus on the text regions. This behavior might not generalize to other types of visualizations and tasks. Click data, gathered via the BubbleView interface, is not uniform over elements, unlike explicit bounding box annotations (i.e., as in the GDI dataset [154]). While this might be a better approximation to natural viewing, non-uniform importance across design elements might cause side-effects for downstream applications like thumbnailing, by cutting off parts of elements or text.

### ■ 9.4  Future directions

This chapter presented the first neural network model for predicting saliency or importance in graphic designs and information visualizations, capable of generalizing to a wide range of design formats. To train this model, we curated hundreds of examples of graphic designs and information visualizations, annotated with importance. The methodology and models can easily be adapted to other visual domains, such as websites [113] or mobile applications [47]. As better webcam-based eyetracking methods become available (e.g., [123, 160, 224]) possibilities also open up for directly training our model from eye movement data.

Chapter 10

# Design applications

I N this chapter we demonstrate how the automatic importance models from the previous chapter can enable diverse design applications. An importance map can provide a common building block for different summarization tasks, including retargeting, thumbnailing, and interactive design tools. The prototypes presented here are meant as proofs-of-concept, showing that our importance prediction alone can give good results with almost no additional post-processing.

## ■ 10.1 Retargeting

The retargeting task is to take a graphic design as input, and to produce a new version of that design with specific dimensions. Retargeting is a common task for modern designers, who must work with many different output dimensions. There is a substantial amount of work on automatic retargeting for natural images, e.g., [6, 186]. Several of these methods have shown that saliency or gaze provide good cues for retargeting, to avoid cropping out image content that people are likely to pay most attention to, such as faces in photographs.

The only previous work on retargeting graphic designs is by O'Donovan et al. [154]. They assumed knowledge of the underlying vector representation of the design and used an expensive optimization with many different energy terms. The method we propose uses bitmap data as input, and is much simpler, without requiring any manual annotations of the input image.

Importance-based retargeting for graphic designs should preserve the most important regions of a design, such as the title and key visual elements. Given a graphic design bitmap as input and specific dimensions, we use the predicted importance map to automatically select a crop of the image with highest importance (Fig. 10.1). Alternative variants of retargeting (e.g., seam carving) are discussed in the Supplemental Material of [35].

**Evaluation:** We ran MTurk experiments where 96 participants were presented with a design and 6 retargeted variants, and were asked to score each variant using a 5-point Likert scale with $1 =$ very poor and $5 =$ very good (Fig. 10.2). Each participant completed this task for 12 designs: 10 randomly selected from a collection of 216 designs, and another 2 designs used for quality control. We used this task to compare crops

**Figure 10.1.** (a) Input designs, (b) our predicted importance maps, and (c) automatic retargeting results using the predicted importance maps to crop out design regions with highest overall importance. This is compared to: (d) edge-based retargeting, where gradient magnitudes are used as the energy map, and (e) Judd saliency, a commonly-used natural image saliency model.

retargeted using predicted importance to crops retargeted using ground truth GDI annotations, Judd saliency [106], DeepGaze saliency [129], and an edge energy map. We extracted a crop with an aspect ratio of 1:4 from a design using the highest-valued region, as assigned by each of the saliency/importance maps. As a baseline, we selected a random crop location.

After an analysis of variance showed a significant effect of retargeting method on score, we performed Bonferroni paired t-tests on the scores of different methods. Across all 216 designs, crops obtained using ground truth GDI annotations had the highest score (Mean: 3.19), followed by DeepGaze (Mean: 2.95) and predicted importance (Mean: 2.92). However, the difference between the latter pair of models was not statistically significant. Edge energy maps (Mean: 2.66) were worse, but not significantly; while Judd saliency (Mean: 2.47) and the random crop baseline (Mean: 2.23) were significantly worse in pairwise comparisons with all the other methods ($p < .01$ for all pairs). Results of additional experimental variants are reported in the Supplemental Material of [35].

Our predicted importance outperforms Judd saliency, a natural saliency model commonly used for comparison [137, 154]. Judd saliency has no notion of text. Predicted importance, trained on less than 1K graphic design images, performs on par with DeepGaze, the currently top-performing neural network-based saliency model [31] which has been trained on 10K natural images, including images with text. Both significantly

**Figure 10.2.** MTurk interface for evaluating retargeting results of predicted importance compared to other baselines.

outperform the edge energy map, which is a common baseline for retargeting. These results show the potential use case of predicted importance for a retargeting task, even without any post-processing steps.

## ■ 10.2 Thumbnailing

Thumbnailing is similar to retargeting, but with a different goal. It aims to provide a visual summary for an image to make it easier to find relevant images in a large collection [104, 210]. Unlike previous methods, our approach operates directly on a bitmap input, rather than requiring a specialized representation as input. For this example our domain is information visualizations rather than graphic designs.

Given an information visualization and an automatically-computed importance map as input, we generate a thumbnail by carving out the less important regions of the image. The importance map is used as an energy function, whereby we iteratively remove image regions with least energy first. Rows and columns of pixels are removed until the desired proportions are achieved, in this case a square thumbnail. This is similar to seam carving [6, 186], but using straight seams, found to work better in our setting. The boundaries of the remaining elements are blurred using the importance map as an alpha-mask with a fade to white. This was done for simplicity, but note that a better design might be to take the highest frequency color as the background. Qualitatively, the resulting thumbnails consist of titles and other main supporting text, as well as data extremes (from the top and bottom of a table, for instance, or from the left and right sides of a plot).

**Evaluation:** We designed a task intended to imitate a search through a database of visualizations. MTurk participants were given a description and a grid of 60 thumbnails, and were instructed to find the visualization that best matches the description. We ran two versions of the study: with the original visualizations resized to thumbnails (Fig. 10.3a), and another with our automatically-computed importance-based thumbnails (Fig. 10.3b). We measured how many clicks it took for participants to find the visualization corresponding to the description in each version.

**Figure 10.3.** Given a set of information visualizations (a), we use our importance maps to automatically generate thumbnails (b). The thumbnails facilitate visual search through a database of visualizations by summarizing the most important content.

A total of 400 participants were recruited for our study. After filtering, we compared the performance of 200 participants who performed the study with resized visualizations and 169 participants who saw the importance-based thumbnails.

Each MTurk assignment, containing a single search task assigned to a single participant, was treated as a repeated observation. We ran an unpaired two-sample t-test to compare the task performance of both groups. On average, participants found the visualization corresponding to the description in fewer clicks using the importance-based thumbnails (1.96 clicks) versus using the resized visualizations (3.25 clicks, t(367) = 5.10, p < .001). Our importance-based thumbnails facilitated speedier retrieval, indicating that the thumbnails captured visualization content relevant for retrieval.

## ■ 10.3 Color theme extraction

*Note: this section was omitted from the UIST'17 paper and appears in this thesis only. Initial user studies performed with this application were not conclusive.*

A color theme is a collection of a small number of colors [137, 153]. Color themes are often used by designers when defining the style of a design. Designers often work from examples; here we propose a method to automatically extract a color theme from a graphic design bitmap. Previous work has shown that, when people create color themes from natural images, they tend to select colors from salient regions [100, 137].

We apply this idea to designs: our method obtains a color theme by sampling the regions predicted most important in a design. In this way, the regions people would most likely pay attention to in a design would contribute most to the color theme representing the design. In Fig. 10.4, we compare a color palette obtained by $k$-means clustering all image pixels ($k = 5$, most commonly) to one obtained by clustering image pixels lying within the top 25% pixels with highest predicted importance. Fine structures in a design, such as text, do not take up many pixels overall, but have high importance, based on importance labels and clicking patterns. It is left to future work to combine

**Figure 10.4.** Some examples where importance can help retrieve more relevant colors for graphic design color themes. We include a k-means baseline, where all image pixels are clustered to generate the colors, compared to natural image saliency and importance-based color themes. For simplicity, for the latter two strategies we perform k-means clustering on the top 25% of the saliency and importance maps, respectively. Importance picks up more on the text in graphic designs, especially the most important text, and ranks the colors accordingly.

importance with other features of a design, and optimize color theme extraction as in the work of Lin and Hanrahan [137].

## ■ 10.4 Interactive applications

An attractive aspect of neural network models is their fast run-time performance. For instance, on a Titan-X GPU, our model computes the importance map for a design in the GDI dataset (600×450 pixels) in 100 ms. Table 10.1 provides some timing information for our model on differently-sized images.

As a prototype, we integrated our importance prediction with a simple design layout tool that allows users to move and resize elements, as well as change color, text font, and opacity (Fig. 9.2). With each change in the design, an importance map is recomputed automatically to provide immediate feedback to the user. The accompanying video and demo (`visimportance.csail.mit.edu`) demonstrate the interactive capabilities of our predictions. The experiments in Chapter 9 (Sec. 9.3.3) provide initial evidence that our model can generalize to the kind of fine-grained design manipulations, like the resizing and relocation of design elements, that would be common in an interactive setting. Determining how best to use importance prediction to provide feedback to users is an interesting problem for future work. For example, importance prediction could help in formulating automatic suggestions for novice users to improve their designs.

| Image size (pixels/side) | 300 | 600 | 900 | 1200 | 1500 |
|---|---|---|---|---|---|
| Avg. compute time (ms) | 46 | 118 | 219 | 367 | 562 |

**Table 10.1.** Time (in milliseconds) taken by our model to compute an importance map for differently-sized images, averaged over 100 trials.

## ■ 10.5 Future directions

This chapter has shown that the same underlying representation of importance, as pixel-wise importance scores, can be used for a variety of design applications, including

retargeting, thumbnailing, extracting color themes, and providing feedback in inter-active design settings. The applications in this chapter are simple prototypes, with limited to no post-processing of the results, and can be improved significantly. Adding an extra segmentation step or else working directly with vector data can allow designs to be more seamlessly edited based on importance. Currently, the importance maps are not guaranteed to be uniform over design elements, and the side-effects of this include cut-off or distorted design regions during thumbnailing and retargeting.

Future work can also explore the use of importance predictions to offer more targeted design feedback and to provide automated suggestions to a user. Our current use of importance within a design tool is to offer feedback about where people are likely to look. However, more active suggestions about where elements can be moved based on importance scores can help optimize designs for effectiveness. In an application where a designer is expected to provide minimal interaction, the designer can specify the intended importance values for a set of design elements, and an automated system can iterate on generating a design under these constraints. There is a wide range of opportunities for integrating importance into design applications, from ones where the designer is in-the-loop to completely automated design systems. The work presented here on predicting importance can open up future doors for using A.I. in the service of facilitating creativity.

# Part V

# Parsing Infographics

# Chapter 11

# Visually29K: A curated dataset of infographics

**A** VAST amount of semantic and design knowledge is encoded in graphic designs, which are created to effectively communicate messages about complex and often abstract topics including 'ways to conserve the environment' and 'understanding the financial crisis'. Graphic designs include clipart [232, 233], comics [98], advertisements [95, 218], diagrams [109, 193], and the infographics that are the focus of this part of the thesis (Fig. 11.1a). Expanding the capabilities of computational algorithms to understand graphic designs can therefore complement natural image understanding by motivating a set of novel research questions with unique challenges. In particular, current techniques trained for natural images do not generalize to the abstract visual elements and diverse styles in graphic designs [95, 98, 218].

In Chapter 12, we tackle the challenge of identifying stand-alone visual elements, which we call 'icons'. To adapt to the stylistic, semantic, and scale variations of icons in graphic designs, which differentiate them from objects in natural images, we propose a synthetic data generation approach. We augment background patches in infographics with a dataset of Internet-scraped icons which we use as training data for an icon proposal mechanism (Fig. 11.1b). In Chapter 13, we use the detected icons along with detected text on infographics for topic prediction and to generate automatic multimodal summaries: outputting the text tags and corresponding icons (visual hashtags) that are most representative of an infographic's topics (Fig. 11.1c).

To make these applications possible, in this chapter we present a curated dataset of 29K infographics called *Visually29K*. The work in this part of the thesis along with the dataset that we plan to release to the research community provides some first steps towards the automated understanding of infographics.

## ■ 11.1 Curating the Visually29K dataset

To facilitate computer vision research on infographics, we assembled the *Visually29K* dataset. We scraped 63K static infographic images from the *Visual.ly* website, a community platform for human-designed visual content. Each infographic is hand categorized, tagged, and described by a designer, making it a rich source of annotated data. We cu-

a) *Visually29K* dataset          b) Synthetic data generation          c) Icon proposals & multi-modal summary

**Figure 11.1.** We make 3 contributions: a) We present *Visually29K*, a curated dataset of infographics (this chapter); b) We generate synthetic data by augmenting Internet-scraped icons onto patches of infographics to train an icon proposal mechanism (Chapter 12); c) We evaluate our automatic icon proposals and present a multi-modal summarization application that takes an infographic and outputs the text tags and visual hashtags that are most representative of the infographic's topics (Chapter 13).

rated this dataset to obtain a representative subset of 28,973 images, ensuring sufficient instances per tag (Table 11.1). The tags associated with images are free-form text, so many of the original tags were either semantically redundant or had too few instances. We cut the original heavy-tailed distribution of 19K tags down to 391 tags with at least 50 exemplars, and by merging redundant tags manually using WordNet [150]. Tags range from concepts which have concrete visual depictions (e.g., *car, cat, baby*) to abstract concepts (e.g., *search engine optimization, foreclosure, revenue*). Metadata for this dataset also includes labels for 26 categories (available for 90% of the infographics), titles (99%) and descriptions (94%), available for future applications. Categories are coarser-grained than tags, and include topics such as *environment, technology, sports*. Each infographic is annotated with a single category but multiple tags.

The infographics in *Visually29K* are very large: up to 5000 pixels per side. Over a third of the infographics are larger than $1000 \times 1500$ pixels. Aspect ratios vary between 5:1 and 1:5. Visual and textual elements occur at a variety of scales, and resizing the images for visual tasks may not be appropriate, given that smaller design elements may be lost.

| Dataset | # of tags | Images per tag | Tags per Image | Aspect ratios |
|---|---|---|---|---|
| 63K (full) | 19469 | min=1 max=3784 mean=7.8 | min=0 max=10 mean=3.7 | min=1:20 max=22:1 |
| 29K (curated) | 391 | min=50 max=2331 mean=151 | min=1 max=9 mean=2.1 | min=1:5 max=5:1 |

**Table 11.1.**  *Visually* dataset statistics. We curated the original 63K infographics available on *Visual.ly* to produce a representative dataset with consistent tags and sufficient instances per tag.



**Figure 11.2.**  User interface for collecting human ground truth to evaluate icon detection and classification. Participants were either asked to annotate all icons on an infographic, or to only annotate icons corresponding to a particular tag (e.g., *fish*).

## ■ 11.2  Human annotations of icons

For a subset of infographics from *Visually29K*, we designed two tasks to collect icon annotations to be used as ground-truth for evaluating computational models (Fig. 11.2). In the first task, we asked participants to annotate all the icons on infographics. In the second task, we asked participants to annotate only the icons corresponding to a particular tag, where the tag comes from the set of tags assigned to the entire infographic. We used the annotations from the first task for evaluating icon proposals (Sec. 12.4), and the second task for evaluating visual hashtags (Sec. 13.2).

**Tag-independent icon annotations:** For 1,400 infographics, we asked participants to "put boxes around any elements that look like icons or pictographs". No further definitions of "icon" were provided. This was a time-consuming task, requiring an average of 15 bounding boxes to be annotated per infographic. A total of 45 participants were recruited, producing a total of 21,288 bounding boxes across all 1,400 infographics. We split these annotated infographics into 400 for validation (training experiments in Sec. 12.2) and 1,000 for testing (reporting final performance in Sec. 12.4).

**Annotation consistency:** Because the interpretation of "icon" may differ across participants, we wanted to measure how consistently humans annotate icons on infographics. For each of 55 infographics, we recruited an additional 5 annotators. Annotations of these participants were compared to the original collected annotations (above).

**Figure 11.3.**    Human agreement in annotating icons is not perfect because different people have different interpretations of "icon". Here we include 3 crops from annotated infographics. In the world map crop we notice three strategies: (i) labeling the entire map as an icon, (ii) labeling individual continents, (iii) labeling the circle graphics superimposed on the map. The set of participant annotations used for evaluation purposes are indicated in red. The other colored boxes depict annotations from additional participants asked to complete the same task for consistency analyses.

We use human consistency as an upper bound on computational models. The scores were averaged across participants and images and are reported in Table 12.1. Human precision and recall are not perfect because different people might disagree about whether a particular visual element (e.g., map, embedded graph, photograph) is an icon. They may also disagree when annotating the boundaries of the icon (Fig. 11.3).

   **Tag-conditional icon annotations:** As ground truth for icon classification, we collected finer-grained annotations by giving participants the same task as before, but asking them to mark bounding boxes around all icons that correspond to a specific text tag (Fig. 11.2). We used 544 infographics along with their associated *Visually29K* text tags, to produce a total of 1,110 separate annotation tasks (each task corresponding to a single image-tag pair). From all these tasks, participants indicated that for 275 (25%) there were no icons on the infographic that corresponded to the text tag. For the remaining 835 image-tag pairs, we collected a total of 7,761 bounding boxes from 45 undergraduate students, averaging 9 bounding boxes per image-tag pair. To compute human consistency for this task as well, for 55 infographics (a total of 172 image-tag pairs) we got an additional 5 annotators. This human upper bound is reported in Table 13.1.

## ■ 11.3  Future directions

The *Visually29K* dataset provides a rich source of training data for future computer vision problems. In the following chapters, we tackle text and icon detection, topic prediction, and summarization. Future applications can additionally make use of the titles, captions, and view statistics that are provided along with the Visually infographics. For instance, meta-data on the number of views and likes can be used to train a classifier of popularity (e.g., [48, 110]) and to learn about the attributes of an infographic that

make it effective (see also Chapter 4). Moreover, while the topic and category labels we have in the dataset can tell us what an infographic is about, we can not directly use this data to infer the intention of the infographic (e.g., whether it has been designed to educate, to convince, to surprise, to dispel bias, to advertise a product, etc.). To be able to make these kinds of higher-level predictions about infographics, we may need to gather additional human annotations (in the style of Hussain et al. [95]).

# Chapter 12

# Synthetically trained icon proposals for parsing infographics

**V**ISUAL elements in infographics have important roles to play in effectively communicating content to the viewer, in a memorable and attention-capturing way (Chapter 4). Motivated by studies of human perception, we aim to design computational algorithms that can parse the visual elements inside infographics - for future summarization, captioning, and information retrieval applications. In this chapter, we tackle the challenge of identifying these stand-alone visual elements, which we call 'icons'.

Instead of (class-specific) icon detection, we instead wish to locate all icon-like elements in an image - i.e., to generate icon proposals. To adapt to the stylistic, semantic, and scale variations of icons in graphic designs (Fig. 12.1), which differentiate them from objects in natural images, we propose a synthetic data generation approach. We augment background patches in infographics with a dataset of Internet-scraped icons which we use as training data for an icon proposal mechanism.

We use the term *icon* to refer to any visual element that has a well-defined closed boundary in space and different appearance from the background (i.e., can be segmented as a stand-alone element). This is inspired by how an object is defined by Alexe et al. [2]. Our approach is more related to *objectness* than to object detection in that we



**Figure 12.1.** Examples of stylistic and semantic variations in scraped icons. a) Visually similar icons scraped for different but semantically related tags *medical, doctor, health, hospital, medicine.* b) Icons with varied styles scraped for the tag *dog.* c) Icons with varied semantic representations for the tag *accident.*

are after class-agnostic object proposals: regions in the image containing icons of any class. While finding icons can be useful for graphic designs more generally, here we train and test icon proposals on our own dataset of infographics.

Training an object detector often requires a large dataset of annotated instances, which is a costly manual effort. We took a different approach, leveraging the fact that infographics are digitally-born to generate synthetic training data: we augmented existing infographics from the *Visually29K* dataset (Chapter 11) with Internet-scraped icons. The advantage of this approach is that we can synthesize any amount of training data by repeatedly sampling new windows from infographics and selecting appropriate patches within the windows to paste new icons into.

## ■ 12.1 Synthetic dataset creation

The use of synthetically generated data to train large CNN models has been gaining popularity, e.g., for learning optical flow [30], action recognition [46], overcoming scattering [190], and object tracking [60]. Simulated environments like video games have been used to collect realistic scene images for semantic segmentation [182]. Our work was inspired by a text recognition system which was trained on a synthetic dataset of images augmented with text [75]. Related to our approach, Dwibedi et al. [51] insert segmented objects into real images to learn to detect natural objects in the wild. We leverage the fact that infographics are digitally-born, so augmenting them with more Internet-scraped design elements is a natural step. Tsutsui and Crandall [213] synthesize compound figures by randomly arranging them on white backgrounds to learn to re-detect them. However, the icons we aim to detect occur on top of complex backgrounds, so we need our synthetic data to capture the visual statistics of in-the-wild infographics.

**Collecting icons:** Starting with the 391 tags in the *Visually29K* dataset, we queried Google with the search terms 'dog icon', 'health icon', etc. for each tag. The search returned a wide range of stylistically and semantically varied icon images (Fig. 12.1). We scraped 250K icons with both transparent and non-transparent backgrounds. Only transparent-background icons were used to augment infographics and train our final icon proposal mechanism, while all 250K icons were used for training an icon classifier (Sec. 13.1). We also present results of training an icon proposal mechanism with icons without transparent backgrounds (Sec. 12.4).

**Augmenting infographics:** To create our synthetic data, we randomly sampled $600 \times 600$px windows from the *Visually29K* infographics. Each window was analyzed for patches of low entropy: measuring the amount of texture in a patch to determine if it is sufficiently empty for icon augmentation. Specifically, from a window, a random patch (with varying location and size) was selected, and Canny edge detection was applied to the patch. The resulting edge values were weighted by a Gaussian window centered on the patch, to give more weight to edges in the center of the patch, and summed to quantify the local entropy, with value ranging from 0 to 1. If the entropy

value was below a predefined threshold, the patch was kept, otherwise it was discarded and a new patch was sampled from the window (Fig. 12.2b). A randomly selected icon from our scraped icon collection was augmented onto each valid patch in a window, for a fixed number of patches per window (Fig. 12.2c). An additional constraint required the icon to meet a set contrast threshold with the patch to ensure it would be visually detectable, or else a new icon would be selected.



**Figure 12.2.**  Synthetic data generation pipeline.  a) Icons with transparent backgrounds scraped from Google. b) Patches selected for augmenting icons, using different approaches. The approach on the left allows more overlap of icons with background elements. The approach on the right is more conservative, selecting appropriate patches to add icons to. c) Infographic windows augmented with the scraped icons.

## ■ 12.2  Effect of synthetic data parameters on model performance

We analyzed how different data augmentation strategies affect the icon proposals on a set of 400 validation infographics containing ground truth annotations of 7,020 bounding boxes. We performed a grid search on 4 augmentation parameters, varying them one at a time: (a) number of icons augmented per window, (b) variation in the size of augmented icons, (c) contrast threshold between the icon and the patch: calculated as difference in color variance between the patch and the icon, and (d) entropy threshold for a patch to be considered valid for augmentation.

We tried 5 settings for the number of icons augmented, from 1 to 16, doubling the number of icons for each experiment. We found no statistically significant differences in the mAP scores of the models trained with these settings.  However, increasing the number of icons augmented per patch increases the time required to generate the synthetic data, since we need to find enough valid image patches to paste icons into. We found that higher scale variation during training helps the model detect icons in infographics, which often occur at different scales. By allowing icons to be augmented at sizes ranging from 30 to 480 pixels per side, we achieved the highest mAP scores. Other settings we tried included capping the maximum icon size at 30, 60, 120, and 240 pixels per side.  Icons larger than 480 pixels per side were not practical with our $600 \times 600$px windows.

We found no significant effects of varying the contrast and entropy thresholds independently, while keeping the other augmentation parameters fixed. However, when we disregard both thresholds and place icons entirely at random in the image windows,

the performance of the trained model degrades significantly (Sec. 12.4). For generating icon proposals on test images, we chose the model with the highest mAP score on the 400 validation images, with 4 icons per window and icon sizes varying from 30 to 240 pixels per side.

## ■ 12.3 Learning to propose icons

We can now use our synthetic data to learn to detect icons. We use the Faster R-CNN network architecture [179], although it is worth noting our training procedure with synthetic icon data can be applied with any architecture. Similar to Dwibedi et al. [51], we are motivated by the fact that Faster R-CNN puts more emphasis on the local visual appearance of an object rather than the global scene layout. As a result, the fact that icons can occur at any location on an infographic is not a problem for the approach.

We adapted Faster R-CNN by making three changes: (i) to handle the large size of infographic images, each image was fed as a cascade of crops, and the detections were aggregated; (ii) we changed the last layer to classify only two categories: any type of icon versus background; (iii) early termination was used during training because the network was found to converge in significantly fewer epochs than the original paper. This could be because detecting generic "iconic" regions depends more on low-level information while category-specific details don't need to be learned.

**Multi-scale detection at test-time:** Infographics in the *Visually29K* data-set are large and contain features at different scales. At test-time, we sampled windows from infographics at 3 different scales to be fed into the network. The first scale spans the entire image. For the two subsequent scales, we sampled 4 and 9 windows, respectively, such that (i) windows at each scale jointly cover the entire image, and (ii) neighboring windows overlap by 10%. Before being fed into the network, every window was rescaled to $600 \times 600$px. The predicted detections from each window were thresholded. Detections across multiple scales were aggregated using non-maximal suppression (NMS) with a value of 0.3. Finally, NMS was applied again to combine smaller detections (often parts of icons) to obtain the final predictions.

**Training details:** We used a total of 10K training instances (windows), where each window was provided with bounding boxes corresponding to the synthetically augmented icons. Faster R-CNN was trained for 30K iterations with a learning rate of $10^{-3}$. Each iteration used a single augmented window to generate a mini-batch of 300 region proposals.

## ■ 12.4 Evaluation of icon proposals

To evaluate our icon proposals, we compare to human annotations of icons on 1,000 test infographics. We report performances using standard detection metrics: precision (*Prec*), recall (*Rec*), F-measure, and mAP. To compute precision and recall, we

threshold IOU at 0.5 (as in the VOC challenge [55]). F-measure is defined as:

$$F_\beta = \frac{(1 + \beta^2)Prec \times Rec}{\beta^2 Prec + Rec}$$

We set $\beta = 0.3$ to weight precision more than recall (a common setting [22]).

**Related methods:** Existing state-of-the-art object detection systems typically rely on object proposal mechanisms [63, 177, 179]. While general literature on detection is aligned with our work, class-agnostic object proposal mechanisms like [2] and [165] are more related. In [2], Alexe et al. present the idea of *objectness*: a metric to quantify how likely it is for an image window to contain an object of *any* class as opposed to background. In [165], authors present a method that generates segmentation masks as class-agnostic object proposals. While related to objectness, our icon proposal mechanism directly proposes object proposals with scores, whereas objectness is a metric that is used to score proposals.

We evaluated whether object proposal mechanisms trained on natural images could generalize to the scale and style variations of icons in infographics. Finding that objectness [2] and [165] failed to detect the icons in our infographics (Table 12.1), we also tried the state-of-the-art object detectors: YOLO9000 [177], SSD [140], and Faster R-CNN [179]. To treat their outputs as object proposals (class-agnostic detections), we report any detection above threshold for any object class that these detectors predict. We also compared to class-agnostic object masks [164]. Default parameters were used for Faster R-CNN and YOLO9000. SSD detections were thresholded at three values (0.01, 0.1, 0.6), and we report results on the best setting (0.01). Re-training Faster R-CNN with our synthetic data (our full model) significantly outperformed networks trained on natural images (Table 12.1).

**Evaluation of synthetic data design choices:** To evaluate the contributions of the main design choices in generating our synthetic data, we ran three additional baselines: (a) augmenting icons in random locations on infographics (instead of finding background patches with low entropy and high contrast with icons), (b) augmenting icons without transparent backgrounds, so that when pasted on an infographic, the augmented icons have clearly-visible boundaries, (c) augmenting icons onto white backgrounds, rather than infographic backgrounds. The last baseline is most similar to the approach in Tsutsui and Crandall [213]. From Table 12.1 we see that all three baselines perform significantly worse than our full model, demonstrating the importance of all three of our design choices: pasting icons with (i) transparent backgrounds onto (ii) appropriate background patches of (iii) in-the-wild infographics. We note that the worst performance among the baselines was when icon proposals were not trained with appropriate backgrounds.

## ■ 12.5 Future directions

In this chapter, we have presented an icon proposal mechanism that, after being trained on synthetic data with icons augmented onto patches of infographics, was able to gen-

| Training data | Model | Prec. | Rec. | $F_{0.3}$ | mAP |
|---|---|---|---|---|---|
| Synthetic with icons | Full model (ours) | 38.8 | 34.3 | 43.2 | 44.2 |
| | Random locations | 27.3 | 12.6 | 26.5 | 29.4 |
| | Non transparent icons | 15.3 | 14.3 | 17.5 | 17.8 |
| | Blank background | 9.3 | 27.6 | 15.1 | 14.5 |
| Natural images | YOLO9000 [177] | 13.6 | 7.1 | 12.6 | 13.7 |
| | Faster R-CNN [179] | 11.0 | 6.0 | 10.2 | 11.4 |
| | SSD [140] | 9.3 | 34.2 | 10.0 | 11.4 |
| | Objectness [2] | 2.9 | 5.6 | 3.1 | 3.0 |
| | Sharpmask [164] | 1.1 | 1.4 | 1.2 | 1.1 |
| | Human upper bound | 63.1 | 64.7 | 61.8 | 66.3 |

**Table 12.1.** Model performance at localizing icons in infographics. The first 4 models were trained with synthetic data containing icons. The next 5 models were trained to detect objects in natural images. The human upper bound is a measure of human consistency on this task. All values are listed as percentages.

eralize to spotting icons in new, in-the-wild infographics. All the evaluations carried out in this chapter pertain to infographics, as we had collected human ground truth annotations of icon locations on a subset of infographics from our *Visually29K* dataset (Chapter 11). It remains up for evaluation whether this icon proposal mechanism generalizes to other multimodal documents, for instance for spotting figures in papers, diagrams in textbooks, or pictures in slideshow presentations. This might require extra training to adapt to the differences in style and format, as well as a collection of human ground truth annotations to validate the results.

Our icon proposal mechanism was specifically trained separately from classification. This is the difference between proposals and detections, where the former are class-agnostic and the latter are class-conditional. This makes our pipeline more easily extendable, since different classifiers can be applied to the proposed icons to accomplish different tasks - e.g., to differentiate abstract from photo-realistic icons or to annotate the proposals using a fixed set of labels for a particular application. For instance, since all the infographics in our *Visually29K* dataset cover a fixed set of 391 topics, we can constrain our icon classifier to predict 1 of 391 classes for each of our icons. We use this approach for generating multimodal summaries of infographics in the next chapter.

# Chapter 13

# Generating multimodal summaries of infographics

**D**ETECTING the textual and visual elements in multimodal documents like infographics can facilitate knowledge retrieval, captioning, and summarization applications. As a first step towards infographic understanding, we propose a multimodal summarization application built upon automatically-detected visual and textual elements (Fig. 13.1). Just as video thumbnails facilitate the sharing, retrieval, and organization of complex media files, our multimodal summaries can be used for effectively capturing a visual digest of complex infographics. Given an infographic as input, our multimodal summary consists of textual and visual hashtags representative of an infographic's topics. We define *visual hashtags* as icons that are most representative of a particular text tag. We evaluate the quality of our multimodal summary by separately testing each component of the pipeline against a set of human annotations.

## ■ 13.1 Approach

**Predicting text tags:** We used Google's Cloud Vision optical character recognition [68] to detect and parse the text from infographics. On average, we extracted 236 words per infographic, of which 170 had *word2vec* representations [69, 149]. The 300-dimensional mean *word2vec* of the bag of extracted words was used as the global feature vector of the text for the infographic. This feature vector was fed into a single-hidden-layer neural network for tag prediction. Since each infographic could have multiple tags, we set this up as a multi-label problem with 391 tags.

   **Classifying icon proposals:** We used the ResNet18 architecture [84] pre-trained on ImageNet, and fine-tuned on icons scraped from Google along with their associated tags (Sec. 12.1). Training was set up as a multi-class problem with 391 tag classes. In addition to the icons with transparent backgrounds, icons with non-transparent backgrounds facilitated the generalization of icon classification to automatically-detected icons.

   **Predicting visual hashtags:** For an input infographic, we predict text tags and generate icon proposals. All the proposals are then fed to the icon classifier to produce a 391-dimensional feature vector of tag probabilities. Then, for each predicted text tag,

**Figure 13.1.** Our computational pipeline for parsing an infographic and computing a multimodal summary. a) The output of our fully-automatic annotation system, running text detection and OCR using Google's Cloud Vision API [68] (semi-transparent green boxes), and our icon detection and classification (red outlines). We trained an icon proposal mechanism with synthetic data to make this system possible. The underlying infographic has been faded to facilitate visualization. b) Our multimodal summarization application uses the detected text and icons on an infographic to produce the text tags and visual hashtags most representative of the infographic's topics.



**Figure 13.2.** Visual hashtags for different concepts. We include 6 different tag classes, sorted by mAP. For each tag class, depicted are the top 4 instances with highest classifier confidence for each tag, constrained to come from different images. Also indicated is the total number (N) of icon proposals per tag class.

we return the icon with the highest probability of belonging to that tag class. Fig. 13.2 contains examples of some visual hashtags: the most confident detections for different tag classes. We demonstrate the automatic output of our system in Fig. 13.3: given an infographic, we predict the text tag and corresponding visual hashtag.

**Input infographics with detections**

Ground truth tag: #happiness   #virus   #cell phone   #oil   #credit card

**Legend:**

☐ (red) Ground truth visual hashtags

☐ (dashed blue) Predicted iconness

☐ (blue) Predicted visual hashtags

**Multi-modal summaries**

Predicted tags: #happiness   #virus   #mobile phone   #economy   #australia

Predicted visual hashtags

#wellness   #mobile   #cell phone   #oil   #credit card

Additional visual hashtags

a)      b)      c)      d)      e)

**Figure 13.3.** Examples of our automated multimodal summarization pipeline, which given an infographic as input, predicts text tags and corresponding visual hashtags. In both (a) and (b), the predicted text tags for the infographics are correct, and the predicted visual hashtags (solid blue boxes) overlap with human annotations (red boxes). Because a single tag might not be sufficient to summarize an infographic, we also provide an additional predicted text tag (second most likely) and corresponding visual hashtag for (a) and (b). In (c)-(e) the text model predicts the wrong tag. In (c), the semantic meaning of the predicted tag is preserved, so the visual hashtag is still correct. In (d) and (e), the wrong visual hashtags are returned as a result of the text predictions. However, we show that if the correct text tag would have been used (bottom, red), correct visual hashtags would have been returned. In dashed blue are all our icon proposals for each infographic. The underlying infographics have been faded to facilitate visualization.

## ■ 13.2 Evaluation

Given an infographic, to evaluate the quality of our predicted text tags, we compared them to the ground truth tags in the *Visually29K* dataset. To evaluate the ability of our computational system to output a relevant visual hashtag for a given infographic and tag, we compare against the human annotations for 544 *Visually29K* infographics (Sec. 11.2). Similar to the task that our computational system receives, participants were asked to annotate all icons corresponding to a particular text tag on an infographic.

**Evaluation of text tag prediction:** Each infographic in *Visually29K* comes with 1-9 tags (2 on average). We achieved 42.6% top-1 average precision and 24.6% top-1 average recall at predicting at least one of an infographic's tags.

**Evaluation of icon classification:** Before evaluating visual hashtags on a per-infographic basis, we evaluate the ability of the icon classifier to retrieve relevant icons *across infographics*. For each of 391 tags, we used the icon classifier's confidence to re-rank all the icon proposals extracted from 544 infographics. Fig. 13.2 contains the highest confidence icon proposals for a few different tags. For each icon proposal, we measure overlap with human annotations: if an icon proposal sufficiently overlaps with a ground truth bounding box (IOU$> 0.5$), that proposal is considered successful. We obtained a mAP of 25.1% by averaging the precision of all the retrieved icon proposals (across all tags).

**Evaluation of visual hashtags:** Next we evaluate the ability of the classifier to retrieve visual hashtags: icons representative of a particular tag, on a *per-infographic* basis. For the following evaluation, we assume that the input is an infographic and a text tag (or multiple tags, if they exist). In a fully automatic setting, the text tags would be predicted by the text model. Fig. 13.3 contains sample results from the fully automatic pipeline. Here we evaluate the quality of our proposed visual hashtags independently of the text model's performance.

For the 835 image-tag pairs with human annotations, we computed the IOU of each of our predicted hashtags with ground truth. We evaluated precision as the percent of predicted visual hashtags that have an IOU $> 0.5$ with at least one of the ground truth annotations. Human participants may annotate multiple icons for an image-tag pair (Fig. 13.3, red boxes). Our application is intended to return a single visual hashtag for a given image-tag pair (Fig. 13.3, solid blue boxes), so we report top-1 precision (Table 13.1). However, we also include the mAP score by considering all our proposals per image-tag pair (Fig. 13.3, dashed blue boxes). From Table 13.1 we see that sorting the icon proposals using our icon classifier produces more relevant results (mAP $= 18.0\%$) for a given tag than just returning the most confident (class-agnostic) icon proposals (mAP $= 14.5\%$). We also verify again that the icon proposals generated by training with icons with *transparent* backgrounds augmented onto *appropriate background regions* of *in-the-wild infographics* outperform baselines that were also trained with synthetic data but with one of these aspects missing.

| Model | Top-1 Prec. | mAP |
|---|:---:|:---:|
| Icon proposals + classification | 27.2 | 18.0 |
| Random locations + class. | 16.7 | 14.2 |
| Non transparent icons + class. | 15.9 | 14.5 |
| Blank background + class. | 16.2 | 14.5 |
| Icon proposals | 16.2 | 14.5 |
| Human upper bound | 55.4 | 57.2 |

**Table 13.1.** Given an infographic and text tag as input, we evaluate the visual hashtags returned. For each image-tag pair, we compute IOU with the ground truth bounding box annotations. A successful visual hashtag is one that has an IOU > 0.5 with at least one ground truth bounding box. All values are listed as percentages.

## ■ 13.3  Future directions

The tools developed in this and the previous chapter provide the ability to densely annotate infographic images - that is, to detect and parse all the text elements, and to detect and classify all the icons. We provide two additional examples of densely annotated infographics in Fig. 13.4. Our current approach considers each text and visual element in isolation, without making use of the strong context afforded by incorporating the semantics of nearby elements. This would help reduce the number of icon misclassifications.

While our main use cases have been infographics, text and icon detection should generalize to other multimodal documents like articles with images and slideshow presentations. Some additional fine-tuning of the training might be required to handle the image types and icon styles more common in articles, for instance to be able to detect larger-scale photographs.

Although we have focused on text and icon annotations as input to a multimodal summarization application, many other applications can be made possible, including automatic captioning and visual question answering. Having an understanding of all the components on an infographic can allow us to generate descriptions for the visually impaired or to facilitate search through a database of infographics. Infographics themselves provide effective summaries on a broad range of topics, so they can be used for information retrieval and building up knowledge databases for artificial agents.
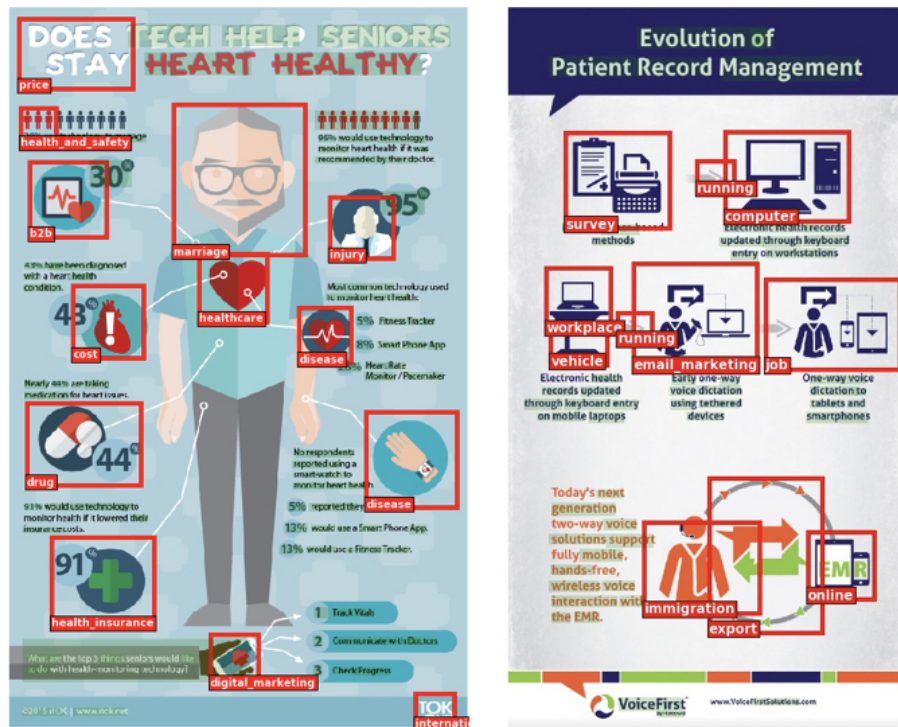
**Figure 13.4.** The output of our fully-automatic dense annotation system, running text detection and OCR using Google's Cloud Vision API [68] (semi-transparent green boxes; we exclude the actual OCR transcriptions here), and our icon detection and classification (red outlines).
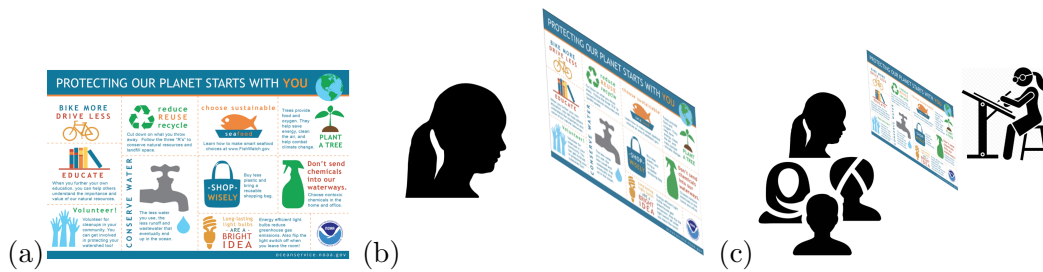
# Part VI

# Conclusion

**Figure 13.5.** In this thesis, we have developed computer vision tools to parse the visual and textual elements inside multimodal documents like infographics by operating directly on the pixels (a). Stepping back, we have also considered the visual design in the context of the viewer (b), by modeling the attention patterns of viewers on information visualizations and graphic designs. Future steps involve stepping back further to examine the individual user context (c), by incorporating the cultural differences, prior experiences, etc. of the user to make individualized predictions. Taking into account the designer that is behind-the-scenes of the design, focus can also turn to the design decisions and intentions that underlie the designs.

**M**ULTIMODAL documents like graphic designs, information visualizations, and infographics are specifically designed with a human viewer in mind, characterized by higher-level semantics, such as a story or a message. Designers choose layouts for the visual and textual elements to effectively convey this story. Developing computational systems for understanding multimodal documents therefore requires tools beyond text and object detection: there needs to be an understanding of the included text (e.g., what topics are discussed, what is the intent), the role of the visual elements (e.g., as concrete examples, to put more emphasis on certain topics, as metaphors, for symbolism), and the relationships between the elements (the layout, the relative sizes and locations of the elements), all of which come together to convey the message.

In this thesis, we scratched the surface of multimodal document understanding, presenting computational models that parse the text in infographics to predict the main topics, detect and classify the visual elements, and put the most representative text and visual elements together into a multimodal summary (Chapters 11-13). This approach considered the image in isolation, making a prediction from pixels to concepts directly using computer vision tools (Fig. 13.5a). We also stepped back from the image to consider it in the context of the viewer (Fig. 13.5b), by modeling the attention patterns of viewers on information visualizations and graphic designs (Chapter 9). We used the predicted attention patterns as input to design applications like automatic feedback within design tools, and automatic summarization in the form of design retargeting and thumbnailing (Chapter 10). Both sets of approaches, for parsing the text and visuals on infographics, and for predicting attention on visualizations and graphic designs, assume no knowledge of the underlying vectorized content. On the one hand, it makes these approaches ambivalent to the encoding format and generalizable to all bitmap images. On the other hand, there is a lot of design information encoded in structured formats like SVG or JavaScript/HTML that can be leveraged for prediction or redesign.

Behind the scenes of every design is a designer who arranged the elements with intention (Fig. 13.5c). Parsing the structured format of a design can make it easier to segment individual elements for further analysis and to evaluate the relative importance of each element. Our current approaches predict importance at the pixel level (Chapter 9), losing sight of semantics. Our current approaches also ignore the relationships between the elements and the layout, which can improve predictions by adding contextual information. Next steps to a deeper understanding of multimodal documents include recognizing metaphors, symbolism, and artistic intent. For instance, consider the infographic in Fig. 13.6 which communicates statistics about microblogging website usage, in the visual format of an iceberg. When we used our computational models to make topic predictions about this infographic (Chapter 13), the visual features were predictive of "travel", while the textual features predicted "social media". Currently, in our approach, we allow the textual predictions to dominate. However, a better description of this infographic would involve recognizing that the iceberg is used as a visual metaphor to make a point about social media. Integrating the visual and textual information on a design can facilitate applications like captioning and visual question answering (e.g., for the visually impaired) and general information retrieval and search (e.g., extracting facts directly from infographics to answer search queries). Human designers are experts at piecing together elements that are cognitively salient or memorable and maximize the utility of information. The space of multimodal documents can give computer vision researchers the opportunity to model and understand the higher-level properties of textual and visual elements in the story being told.

In this thesis, the approach to image saliency and importance of visual content has focused on modeling patterns averaged across a population (Chapters 6-9), without stepping back to consider the individual user context, including cultural differences and prior experience (Fig. 13.5c). The human-computer interfaces we designed were used to collect attention and importance data from many participants (Chapters 7-8), and the average attention maps were used for analyses, to train automatic models, and to re-target and thumbnail designs (Chapters 9-10). Building interfaces to crowdsource user attention data at-scale can facilitate more individualized approaches: sufficient data per user can be used for generating multimodal documents and document summaries that are customized to the user. For instance, knowing where a particular individual would look or what information they would find important can be used to more effectively design content for the user (a sample application for generating personalized GIFs of academic posters is prototyped in Chapter 8). This thesis has focused mainly on discriminative problems: making inferences about multimodal documents taken as input. Future work can leverage the computational understanding built up about multimodal documents for generative problems: synthesizing multimodal content like diagrams, posters, and graphs from raw text and data. Other ideas include translating textbooks into interactive slideshow presentations and making articles more interactive by populating them with relevant images and visualizations (see Chapter 1).

The space of multimodal documents has traditionally been tackled by separate re-

**Figure 13.6.** This infographic (courtesy of *Digital Jungle*) communicates statistics about microblogging website usage in the visual format of an iceberg. When we used our computational models to make topic predictions about this infographic, the visual features were predictive of "travel", while the textual features predicted "social media". Currently, in our approach, we allow the textual predictions to dominate. However, a better description of this infographic would involve recognizing that the iceberg is used as a visual metaphor to make a point about social media. Integrating the visual and textual information on a design can facilitate applications like captioning and visual question answering.

search communities, working on text parsing, image recognition, and style/format classification (Chapter 2). In this thesis, tools and ideas from computer vision, human perception, natural language processing, and human-computer interaction have been combined and adapted for parsing information visualizations, graphic designs, and infographics. The design and summarization applications demonstrated in this thesis are just the tip of the iceberg for what can be done when the visual and textual elements on multimodal documents can be fully segmented and understood. This thesis has laid some of the groundwork and tools necessary to make more applications possible in the future, for instance for democratizing graphic design for internet-scale education and for disseminating information in a broadly-accessible, multimodal format.

## ■ 13.4 Contributions

This thesis has presented contributions at the interface of human and computer vision, with applications to human-computer interfaces and design. Specifically, in this thesis, we have presented:

- An analysis of where people look most and what they find interesting in natural images (Chapters 6, 8.3), information visualizations (Chapters 4, 9.3.1), and graphic designs (Chapter 9.3.2). These findings can be used to design automated models for predicting attention and importance on images.

- The notion of "importance" to replace saliency as a predictor of image regions people are likely to look at over longer periods of time and find interesting (Chapters 1.3, 6, 7).

- Two novel data collection methods for crowdsourcing importance based on a moving-window methodology, using mouse clicks (Chapter 7) and zoom gestures (Chapter 8). Both can be used for efficiently scaling up the collection of attention data online.

- A neural network model to predict importance on novel graphic designs and information visualizations (Chapter 9), capable of being run in real-time within interactive design tools (Chapter 10.4).

- Automatic pipelines for retargeting, thumbnailing, and otherwise summarizing posters (Chapter 8.4), graphic designs (Chapter 10.1), information visualizations (Chapter 10.2), and infographics (Chapter 13).

- A large-scale curated dataset of infographics, with associated tags and human annotations, to facilitate computer vision research on non-natural images (Chapter 11).

- A synthetically trained neural network model for detecting icons in infographics (Chapter 12).

- A computational pipeline for parsing the textual and visual elements in infographics in order to predict the main topics of the infographics (Chapter 13).

## ■ 13.5  Decade-long vision

Developing the computational tools to parse a wide range of multimodal documents will allow A.I. agents to query knowledge from and answer questions about a more diverse set of topics than found in text documents and photo collections alone. Computational understanding precedes synthesis. I envision that the next generation of algorithms will be able to automatically create novel multimodal content, combining not just visuals and text, but sounds and motions, physical and virtual experiences, to explain concepts, imagine and tell stories. These algorithms could then harness knowledge of existing designs and automatically-learned metaphors and analogies, to visualize abstract concepts not previously visualized by humans. Building up knowledge about individual users could be used to build systems that synthesize personalized content, efficiently educating and delivering information to diverse audiences. Currently, designing audio-visual and interactive content to communicate complex concepts is done manually by trained designers. Augmenting designers with automatic systems will allow for the greater democratization of information. Being able to take any piece of information and automatically deliver it to any user in personalized physical or virtual reality form will have important educational implications.

# Bibliography

[1] Amer Al-Rahayfeh and Miad Faezipour. Eye tracking and head movement detection: A state-of-art survey. *IEEE Journal of Translational Engineering in Health and Medicine*, 1, 2013. ISSN 2168-2372. doi: 10.1109/JTEHM.2013.2289879.

[2] Bogdan Alexe, Thomas Deselaers, and Vittorio Ferrari. Measuring the objectness of image windows. *IEEE TPAMI*, 34(11):2189–2202, 2012.

[3] Basak E Alper, Nathalie Henry Riche, and Tobias Hollerer. Structuring the space: a study on enriching node-link diagrams with visual references. In *CHI '14*, pages 1825–1834. ACM, 2014.

[4] Christine Alvarado and Randall Davis. Sketchread: a multi-domain sketch recognition engine. In *UIST*, pages 23–32. ACM, 2004.

[5] Apostolos Antonacopoulos, David Bridson, Christos Papadopoulos, and Stefan Pletschacher. A realistic dataset for performance evaluation of document layout analysis. In *ICDAR*, pages 296–300. IEEE, 2009.

[6] Shai Avidan and Ariel Shamir. Seam carving for content-aware image resizing. *ACM Trans. Graph.*, 26(3), July 2007. ISSN 0730-0301. doi: 10.1145/1276377. 1276390. URL http://doi.acm.org/10.1145/1276377.1276390.

[7] Shumeet Baluja and Dean Pomerleau. Non-intrusive gaze tracking using artificial neural networks. In *Advances in Neural Information Processing Systems*, pages 753–760, 1994.

[8] Scott Bateman, Regan L Mandryk, Carl Gutwin, Aaron Genest, David McDine, and Christopher Brooks. Useful junk?: the effects of visual embellishment on comprehension and memorability of charts. In *CHI '10*, pages 2573–2582, 2010.

[9] Roman Bednarik and Markku Tukiainen. Effects of display blurring on the behavior of novices and experts during program debugging. In *CHI'05 EA*, pages 1204–1207. ACM, 2005.

[10] Roman Bednarik and Markku Tukiainen. Validating the restricted focus viewer: A study using eye-movement tracking. *Behavior research methods*, 39(2):274–282, 2007.

[11] Jennifer Romano Bergstrom and Andrew Schall. *Eye tracking in user experience design*. Elsevier, 2014.

[12] Kathrin Berkner, Edward L Schwartz, and Christophe Marle. Smartnails: display- and image-dependent thumbnails. In *Document Recognition and Retrieval XI*, volume 5296, pages 54–66. International Society for Optics and Photonics, 2003.

[13] Alan F Blackwell and TRG Green. Does metaphor increase visual language usability? In *Visual Languages*, pages 246–253. IEEE, 1999.

[14] Alan F. Blackwell, Anthony R. Jansen, and Kim Marriott. *Restricted Focus Viewer: A Tool for Tracking Visual Attention*, pages 162–177. Springer Berlin Heidelberg, Berlin, Heidelberg, 2000. ISBN 978-3-540-44590-6. doi: 10.1007/ 3-540-44590-0_17. URL http://dx.doi.org/10.1007/3-540-44590-0_17.

[15] T Blascheck, K Kurzhals, M Raschke, M Burch, D Weiskopf, and T Ertl. State-of-the-art of visualization for eye tracking data. In *Proceedings of EuroVis*, volume 2014, 2014.

[16] Rita Borgo, Alfie Abdul-Rahman, Farhan Mohamed, Philip W Grant, Irene Reppa, Luciano Floridi, and Min Chen. An empirical study on using visual embellishments in visualization. *IEEE TVCG*, 18(12):2759–2768, 2012.

[17] A. Borji and L. Itti. State-of-the-art in visual attention modeling. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 35(1):185–207, 2013.

[18] A. Borji, D. N. Sihite, and L. Itti. Quantitative analysis of human-model agreement in visual saliency modeling: a comparative study. *IEEE Transactions on Image Processing*, 22(1):55–69, 2013.

[19] Ali Borji. Boosting bottom-up and top-down visual features for saliency estimation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 438–445, 2012.

[20] Ali Borji and Laurent Itti. Cat2000: A large scale fixation dataset for boosting saliency research. *arXiv preprint arXiv:1505.03581*, 2015.

[21] Ali Borji, Hamed R. Tavakoli, Dicky N. Sihite, and Laurent Itti. Analysis of scores, datasets, and models in visual saliency prediction. In *International Conference on Computer Vision*, December 2013.

[22] Ali Borji, Ming-Ming Cheng, Huaizu Jiang, and Jia Li. Salient object detection: A benchmark. *IEEE TIP*, 24(12):5706–5722, 2015.

[23] Michelle A. Borkin, Azalea A Vo, Zoya Bylinskii, Phillip Isola, Shashank Sunkavalli, Aude Oliva, and Hanspeter Pfister. What makes a visualization memorable? *IEEE Transactions on Visualization and Computer Graphics*, 19(12): 2306–2315, Dec 2013. ISSN 1077-2626. doi: 10.1109/TVCG.2013.234.

[24] Michelle A. Borkin, Zoya Bylinskii, Nam Wook Kim, Constance May Bainbridge, Chelsea S. Yeh, Daniel Borkin, Hanspeter Pfister, and Aude Oliva. Beyond memorability: Visualization recognition and recall. *IEEE Transactions on Visualization and Computer Graphics*, 22(1):519–528, Jan 2016. ISSN 1077-2626. doi: 10.1109/TVCG.2015.2467732.

[25] Léon Bottou. Stochastic learning. In *Advanced lectures on machine learning*, pages 146–168. Springer, 2004.

[26] Neil Bruce and John Tsotsos. Attention based on information maximization. *Journal of Vision*, 7(9):950–950, 2007.

[27] Daniel Bruneau, M Angela Sasse, and JD McCarthy. The eyes never lie: The use of eye tracking data in hci research. In *Proceedings of the CHI*, volume 2, page 25, 2002.

[28] M. Burch, N. Konevtsova, J. Heinrich, M. Hoeferlin, and D. Weiskopf. Evaluation of traditional, orthogonal, and radial tree diagrams by an eye tracking study. *IEEE TVCG*, 17(12):2440–2448, Dec 2011. ISSN 1077-2626. doi: 10.1109/TVCG.2011. 193.

[29] Georg Buscher, Edward Cutrell, and Meredith Ringel Morris. What do you see when you're surfing?: Using eye tracking to predict salient regions of web pages. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, CHI '09, pages 21–30, New York, NY, USA, 2009. ACM. ISBN 978-1-60558-246-7. doi: 10.1145/1518701.1518705. URL http://doi.acm.org/10.1145/1518701. 1518705.

[30] D. J. Butler, J. Wulff, G. B. Stanley, and M. J. Black. A naturalistic open source movie for optical flow evaluation. In *ECCV*, 2012.

[31] Zoya Bylinskii, Tilke Judd, Ali Borji, Laurent Itti, Frédo Durand, Aude Oliva, and Antonio Torralba. MIT Saliency Benchmark, 2012.

[32] Zoya Bylinskii, Ellen M. DeGennaro, Rishi Rajalingham, Harald Ruda, Jinxia Zhang, and John K. Tsotsos. Towards the quantitative evaluation of visual attention models. *Vision Research*, 116, Part B:258 – 268, 2015. ISSN 0042-6989. doi: https://doi.org/10.1016/j.visres.2015.04.007. URL http://www.sciencedirect. com/science/article/pii/S0042698915001522.

[33] Zoya Bylinskii, Adrià Recasens, Ali Borji, Aude Oliva, Antonio Torralba, and Frédo Durand. Where should saliency models look next? In *European Conference on Computer Vision*, pages 809–824. Springer, 2016.

[34] Zoya Bylinskii, Michelle A. Borkin, Nam Wook Kim, Hanspeter Pfister, and Aude Oliva. Eye fixation metrics for large scale evaluation and comparison of information visualizations. In Michael Burch, Lewis Chuang, Brian Fisher, Albrecht Schmidt, and Daniel Weiskopf, editors, *Eye Tracking and Visualization: Foundations, Techniques, and Applications. ETVIS 2015*, pages 235–255. Springer International Publishing, Cham, 2017. doi: 10.1007/978-3-319-47024-5_14.

[35] Zoya Bylinskii, Nam Wook Kim, Peter O'Donovan, Sami Alsheikh, Spandan Madan, Hanspeter Pfister, Fredo Durand, Bryan Russell, and Aaron Hertzmann. Learning visual importance for graphic designs and data visualizations. In *Proceedings of the 30th Annual ACM Symposium on User Interface Software & Technology*, UIST '17. ACM, 2017. ISBN 978-1-4503-4981-9/17/10. doi: 10.1145/3126594.3126653. URL https://doi.org/10.1145/3126594.3126653.

[36] Zoya Bylinskii, Tilke Judd, Aude Oliva, Antonio Torralba, and Frédo Durand. What do different evaluation metrics tell us about saliency models? *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2018. doi: 10.1109.

[37] Deng Cai, Shipeng Yu, Ji-Rong Wen, and Wei-Ying Ma. Vips: a vision-based page segmentation algorithm. *Microsoft Research Technical Report*, 2003.

[38] Moran Cerf, E Paxon Frady, and Christof Koch. Faces and text attract gaze independent of the task: Experimental data and computer model. *Journal of vision*, 9(12), 2009.

[39] Mon Chu Chen, John R. Anderson, and Myeong Ho Sohn. What can a mouse cursor tell us more?: Correlation of eye/mouse movements on web browsing. In *CHI '01 Extended Abstracts on Human Factors in Computing Systems*, CHI EA '01, pages 281–282, New York, NY, USA, 2001. ACM. ISBN 1-58113-340-5. doi: 10.1145/634067.634234. URL http://doi.acm.org/10.1145/634067.634234.

[40] William S Cleveland and Robert McGill. Graphical perception: Theory, experimentation, and application to the development of graphical methods. *J. of the Am. Stat. Assoc.*, 79(387):531–554, 1984.

[41] Marcella Cornia, Lorenzo Baraldi, Giuseppe Serra, and Rita Cucchiara. A deep multi-level network for saliency prediction. In *2016 23rd International Conference on Pattern Recognition (ICPR)*, pages 3488–3493, Dec 2016. doi: 10.1109/ICPR. 2016.7900174.

[42] Laura Cowen, Linden Js Ball, and Judy Delin. An eye movement analysis of web page usability. In *People and Computers XVI*, pages 317–335. Springer, 2002.

[43] Edward Cutrell and Zhiwei Guan. What are you looking for?: An eye-tracking study of information usage in web search. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, CHI '07, pages 407–416, New York, NY, USA, 2007. ACM. ISBN 978-1-59593-593-9. doi: 10.1145/1240624. 1240690. URL http://doi.acm.org/10.1145/1240624.1240690.

[44] G.B. Dantzig. Application of the simplex method to a transportation problem. *Activity Analysis of Production and Allocation*, pages 359–373, 1951.

[45] Abhishek Das, Harsh Agrawal, Lawrence Zitnick, Devi Parikh, and Dhruv Batra. Human attention in visual question answering: Do humans and deep networks look at the same regions? *arXiv preprint arXiv:1606.03556*, 2016.

[46] CR De Souza, A Gaidon, Y Cabon, and AM Lopez Pena. Procedural generation of videos to train deep action recognition networks. In *CVPR*, 2017.

[47] Biplab Deka, Zifeng Huang, Chad Franzen, Joshua Hibschman, Daniel Afergan, Yang Li, Jeffrey Nichols, and Ranjitha Kumar. Rico: A mobile app dataset for building data-driven design applications. In *Proceedings of the 30th Annual Symposium on User Interface Software and Technology*, UIST '17, 2017.

[48] Arturo Deza and Devi Parikh. Understanding image virality. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1818–1826, 2015.

[49] Andrew T Duchowski. A breadth-first survey of eye-tracking applications. *Behavior Research Methods, Instruments, & Computers*, 34(4):455–470, 2002.

[50] Andrew T Duchowski. Eye tracking methodology. *Theory and practice*, 328, 2007.

[51] Debidatta Dwibedi, Ishan Misra, and Martial Hebert. Cut, paste and learn: Surprisingly easy synthesis for instance detection. *ICCV*, 2017.

[52] W. Einhäuser and P. Konig. Does luminance-contrast contribute to a saliency for overt visual attention? *European Journal of Neuroscience*, 17:1089–1097, 2003.

[53] Kevin Ellis, Daniel Ritchie, Armando Solar-Lezama, and Joshua B Tenenbaum. Learning to infer graphics programs from hand-drawn images. *arXiv preprint arXiv:1707.09627*, 2017.

[54] U. Engelke, H. Liu, J. Wang, P. Le Callet, I. Heynderickx, H-J Zepernick, and A. Maeder. Comparative study of fixation density maps. *IEEE TIP*, 22(3):1121–1133, 2013.

[55] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The Pascal Visual Object Classes (VOC) Challenge. *IJCV*, 88(2):303–338, 2010.

[56] A. Farhadi, A. Gupta, A. Kembhavi, R. Mottagi, E. Kolve, G. Sigurdsson, J. Choi, D. Schwenk, and D. Gordon. Workshop on visual understanding across modalities. In *CVPR Workshops*, 2017. URL http://vuchallenge.org/index.html.

[57] T. Fawcett. An introduction to ROC analysis. *Pattern recognition letters*, 27(8): 861–874, 2006.

[58] Simone Frintrop, Erich Rome, and Henrik I. Christensen. Computational visual attention systems and their cognitive foundations: A survey. *ACM Trans. Appl. Percept.*, 7(1):6:1–6:39, January 2010. ISSN 1544-3558. doi: 10.1145/1658349. 1658355. URL http://doi.acm.org/10.1145/1658349.1658355.

[59] Kenneth Alberto Funes Mora, Florent Monay, and Jean-Marc Odobez. Eyediap: A database for the development and evaluation of gaze estimation algorithms from rgb and rgb-d cameras. In *Proceedings of the Symposium on Eye Tracking Research and Applications*, ETRA '14, pages 255–258, New York, NY, USA, 2014. ACM. ISBN 978-1-4503-2751-0. doi: 10.1145/2578153.2578190. URL http://doi.acm.org/10.1145/2578153.2578190.

[60] Adrien Gaidon, Qiao Wang, Yohann Cabon, and Eleonora Vig. Virtual worlds as proxy for multi-object tracking analysis. In *CVPR*, 2016.

[61] Elena Garces, Aseem Agarwala, Diego Gutierrez, and Aaron Hertzmann. A similarity measure for illustration style. *ACM Transactions on Graphics (TOG)*, 33 (4):93, 2014.

[62] Sohaib Ghani and Niklas Elmqvist. Improving revisitation in graphs through static spatial features. In *Proceedings of Graphics Interface 2011*, pages 175–182. Canadian Human-Computer Communications Society, 2011.

[63] Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *CVPR*, 2014.

[64] Stas Goferman, Lihi Zelnik-Manor, and Ayellet Tal. Context-aware saliency detection. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 34 (10):1915–1926, Oct 2012. ISSN 0162-8828. doi: 10.1109/TPAMI.2011.272.

[65] Joseph H. Goldberg and Jonathan I. Helfman. Comparing information graphics: A critical look at eye tracking. In *BELIV'10*, pages 71–78. ACM, 2010. ISBN 978-1-4503-0007-0. doi: 10.1145/2110192.2110203. URL http://doi.acm.org.ezp-prod1.hul.harvard.edu/10.1145/2110192.2110203.

[66] Joseph H Goldberg and Xerxes P Kotval. Computer interface evaluation using eye movements: methods and constructs. *International Journal of Industrial Ergonomics*, 24(6):631–645, 1999.

[67] Joseph H. Goldberg, Mark J. Stimson, Marion Lewenstein, Neil Scott, and Anna M. Wichansky. Eye tracking in web search tasks: Design implications. In *Proceedings of the 2002 Symposium on Eye Tracking Research & Applications*, ETRA '02, pages 51–58, New York, NY, USA, 2002. ACM. ISBN 1-58113-467-3. doi: 10.1145/507072.507082. URL http://doi.acm.org/10.1145/507072.507082.

[68] Google. Cloud Vision API: Optical Character Recogition. https://cloud.google.com/vision/, accessed in October 2017.

[69] Google. Word2vec model. https://drive.google.com/file/d/0B7XkCwpI5KDYNlNUTTlSS21pQmM/edit?usp=sharing, accessed in October 2017.

[70] W Graf and H Krueger. Ergonomic evaluation of user-interfaces by means of eye-movement data. In *Proceedings of the third international conference on human-computer interaction*, pages 659–665. Elsevier Science Inc., 1989.

[71] Elizabeth R. Grant and Michael J. Spivey. Eye movements and problem solving. *Psychological Science*, 14(5):462–466, 2003. doi: 10.1111/1467-9280.02454. URL http://dx.doi.org/10.1111/1467-9280.02454.

[72] D. M Green and J. A Swets. *Signal detection theory and psychophysics*. John Wiley, 1966.

[73] Rebecca Grier, Philip Kortum, and James Miller. How users view web pages: An exploration of cognitive and perceptual mechanisms. *Human computer interaction research in Web design and evaluation*, pages 22–41, 2007.

[74] Qi Guo and Eugene Agichtein. Towards predicting web searcher gaze position from mouse movements. In *CHI '10 Extended Abstracts on Human Factors in Computing Systems*, CHI EA '10, pages 3601–3606, New York, NY, USA, 2010. ACM. ISBN 978-1-60558-930-5. doi: 10.1145/1753846.1754025. URL http://doi.acm.org/10.1145/1753846.1754025.

[75] A. Gupta, A. Vedaldi, and A. Zisserman. Synthetic data for text localisation in natural images. In *CVPR*, 2016.

[76] David Ha and Douglas Eck. A neural representation of sketch drawings. *arXiv preprint arXiv:1704.03477*, 2017.

[77] Michael J. Haass, Andrew T. Wilson, Laura E. Matzen, and Kristin M. Divis. *Modeling Human Comprehension of Data Visualizations*, pages 125–134. Springer International Publishing (VAMR), Cham, 2016. ISBN 978-3-319-39907-2. doi: 10.1007/978-3-319-39907-2_12. URL http://dx.doi.org/10.1007/978-3-319-39907-2_12.

[78] Mark A Hall. *Correlation-based feature selection for machine learning.* PhD thesis, The University of Waikato, 1999.

[79] Adam W Harley, Alex Ufkes, and Konstantinos G Derpanis. Evaluation of deep convolutional nets for document image classification and retrieval. In *ICDAR*, pages 991–995. IEEE, 2015.

[80] Steve Haroz and David Whitney. How capacity limits of attention influence information visualization effectiveness. *IEEE TVCG*, 18(12):2402–2410, 2012.

[81] Steve Haroz, Robert Kosara, and Steven L Franconeri. Isotype visualization– working memory, performance, and engagement with pictographs. In *CHI '15*, pages 1191–1200. ACM, 2015.

[82] Lane Harrison, Katharina Reinecke, and Remco Chang. Infographic aesthetics: Designing for the first impression. In *Proceedings of the 33rd Annual ACM Conference on Human Factors in Computing Systems*, CHI '15, pages 1187–1190, New York, NY, USA, 2015. ACM. ISBN 978-1-4503-3145-6. doi: 10.1145/2702123. 2702545. URL http://doi.acm.org/10.1145/2702123.2702545.

[83] Mary Hayhoe. Advances in relating eye movements and cognition. *Infancy*, 6(2): 267–274, 2004.

[84] K. He, X. Zhang, S. Ren, and J. Sun. Deep Residual Learning for Image Recognition. In *CVPR*, 2016.

[85] Kenneth Holmqvist, Marcus Nyström, Richard Andersson, Richard Dewhurst, Halszka Jarodzka, and Joost Van de Weijer. *Eye tracking: A comprehensive guide to methods and measures.* OUP Oxford, 2011.

[86] Xiaodi Hou and Liqing Zhang. Saliency detection: A spectral residual approach. In *2007 IEEE Conference on Computer Vision and Pattern Recognition*, pages 1–8, June 2007. doi: 10.1109/CVPR.2007.383267.

[87] Jeff Huang, Ryen W. White, and Susan Dumais. No clicks, no problem: Using cursor movements to understand and improve search. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, CHI '11, pages 1225–1234, New York, NY, USA, 2011. ACM. ISBN 978-1-4503-0228-9. doi: 10. 1145/1978942.1979125. URL http://doi.acm.org/10.1145/1978942.1979125.

[88] Jeff Huang, Ryen White, and Georg Buscher. User see, user point: Gaze and cursor alignment in web search. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, CHI '12, pages 1341–1350, New York, NY, USA, 2012. ACM. ISBN 978-1-4503-1015-4. doi: 10.1145/2207676.2208591. URL http://doi.acm.org/10.1145/2207676.2208591.

[89] Qiong Huang, Ashok Veeraraghavan, and Ashutosh Sabharwal. Tabletgaze: unconstrained appearance-based gaze estimation in mobile tablets. *arXiv preprint arXiv:1508.01244*, 2015.

[90] Weidong Huang. Using eye tracking to investigate graph layout effects. In *APVIS '07*, pages 97–100, Feb 2007. doi: 10.1109/APVIS.2007.329282.

[91] Weidong Huang, P. Eades, and Seok-Hee Hong. A graph reading behavior: Geodesic-path tendency. In *PacificVis '09*, pages 137–144, April 2009. doi: 10.1109/PACIFICVIS.2009.4906848.

[92] Weihua Huang and Chew Lim Tan. A system for understanding imaged infographics and its applications. In *Proceedings of the 2007 ACM symposium on Document engineering*, pages 9–18. ACM, 2007.

[93] Xun Huang, Chengyao Shen, Xavier Boix, and Qi Zhao. Salicon: Reducing the semantic gap in saliency prediction by adapting deep neural networks. In *2015 IEEE International Conference on Computer Vision (ICCV)*, pages 262–270, Dec 2015. doi: 10.1109/ICCV.2015.38.

[94] Jessica Hullman, Eytan Adar, and Priti Shah. Benefitting infovis with visual difficulties. *IEEE TVCG*, 17(12):2213–2222, 2011.

[95] Zaeem Hussain, Mingda Zhang, Xiaozhong Zhang, Keren Ye, Christopher Thomas, Zuha Agha, Nathan Ong, and Adriana Kovashka. Automatic understanding of image and video advertisements. In *CVPR*, 2017.

[96] Laurent Itti and Christof Koch. Computational modelling of visual attention. *Nature reviews neuroscience*, 2(3):194, 2001.

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

[98] M. Iyyer, V. Manjunatha, A. Guha, Y. Vyas, J. Boyd-Graber, H. Daumé, III, and L. Davis. The Amazing Mysteries of the Gutter: Drawing Inferences Between Panels in Comic Book Narratives. In *CVPR*, 2017.

[99] Robert Jacob and Keith S Karn. Eye tracking in human-computer interaction and usability research: Ready to deliver the promises. *Mind*, 2(3):4, 2003.

[100] Ali Jahanian, S. V. N. Vishwanathan, and Jan P. Allebach. Autonomous color theme extraction from images using saliency. *Proc. SPIE*, 9408:940807–940807–8, 2015. doi: 10.1117/12.2084051. URL http://dx.doi.org/10.1117/12.2084051.

[101] Anthony R Jansen, Alan F Blackwell, and Kim Marriott. A tool for tracking visual attention: The restricted focus viewer. *Behavior Research Methods*, 35(1): 57–69, 2003.

[102] Yangqing Jia, Evan Shelhamer, Jeff Donahue, Sergey Karayev, Jonathan Long, Ross B. Girshick, Sergio Guadarrama, and Trevor Darrell. Caffe: Convolutional architecture for fast feature embedding. *CoRR*, abs/1408.5093, 2014. URL http://arxiv.org/abs/1408.5093.

[103] Ming Jiang, Shengsheng Huang, Juanyong Duan, and Qi Zhao. Salicon: Saliency in context. In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1072–1080, June 2015. doi: 10.1109/CVPR.2015.7298710.

[104] Binxing Jiao, Linjun Yang, Jizheng Xu, and Feng Wu. Visual summarization of web pages. In *Proceedings of the 33rd International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR '10, pages 499–506, New York, NY, USA, 2010. ACM. ISBN 978-1-4503-0153-4. doi: 10.1145/1835449.1835533. URL http://doi.acm.org/10.1145/1835449.1835533.

[105] Sheree Josephson and Michael E. Holmes. Visual attention to repeated internet images: Testing the scanpath theory on the world wide web. In *Proceedings of the 2002 Symposium on Eye Tracking Research & Applications*, ETRA '02, pages 43–49, New York, NY, USA, 2002. ACM. ISBN 1-58113-467-3. doi: 10.1145/507072.507081. URL http://doi.acm.org/10.1145/507072.507081.

[106] Tilke Judd, Krista Ehinger, Frédo Durand, and Antonio Torralba. Learning to predict where humans look. In *2009 IEEE 12th International Conference on Computer Vision*, pages 2106–2113, Sept 2009. doi: 10.1109/ICCV.2009.5459462.

[107] Tilke Judd, Frédo Durand, and Antonio Torralba. A benchmark of computational models of saliency to predict human fixations. In *MIT Technical Report*, 2012.

[108] Marcel Adam Just and Patricia A Carpenter. Eye fixations and cognitive processes. *Cognitive psychology*, 8(4):441–480, 1976.

[109] Aniruddha Kembhavi, Mike Salvato, Eric Kolve, Minjoon Seo, Hannaneh Hajishirzi, and Ali Farhadi. A digram is worth a dozen images. In *ECCV*, 2016.

[110] Aditya Khosla, Atish Das Sarma, and Raffay Hamid. What makes an image popular? In *Proceedings of the 23rd international conference on World wide web*, pages 867–876. ACM, 2014.

[111] Wolf Kienzle, Felix A Wichmann, Matthias O Franz, and Bernhard Schölkopf. A nonparametric approach to bottom-up visual saliency. In *Advances in Neural Information Processing Systems*, pages 689–696, 2006.

[112] Nam Wook Kim, Zoya Bylinskii, Michelle A. Borkin, Aude Oliva, Krzysztof Z. Gajos, and Hanspeter Pfister. A crowdsourced alternative to eye-tracking for visualization understanding. In *Proceedings of the 33rd Annual ACM Conference Extended Abstracts on Human Factors in Computing Systems*, CHI EA '15, pages

1349–1354, New York, NY, USA, 2015. ACM. ISBN 978-1-4503-3146-3. doi: 10. 1145/2702613.2732934. URL http://doi.acm.org/10.1145/2702613.2732934.

[113] Nam Wook Kim, Zoya Bylinskii, Michelle A. Borkin, Krzysztof Z. Gajos, Aude Oliva, Frédo Durand, and Hanspeter Pfister. Bubbleview: an interface for crowd-sourcing image importance maps and tracking visual attention. *TOCHI*, 2017. doi: 10.1145/3131275. URL https://arxiv.org/abs/1702.05150.

[114] Nam Wook Kim, Zoya Bylinskii, Michelle A Borkin, Krzysztof Z Gajos, Aude Oliva, Frédo Durand, and Hanspeter Pfister. Bubbleview: an interface for crowd-sourcing image importance maps and tracking visual attention. *TOCHI*, 2017.

[115] Sung-Hee Kim, Zhihua Dong, Hanjun Xian, Benjavan Upatising, and Ji Soo Yi. Does an eye tracker tell the truth about visualizations?: Findings while investi-gating visualizations for decision making. *IEEE TVCG*, 18(12):2421–2430, 2012. ISSN 1077-2626. doi: 10.1109/TVCG.2012.215.

[116] Aniket Kittur, Ed H. Chi, and Bongwon Suh. Crowdsourcing user studies with mechanical turk. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, CHI '08, pages 453–456, New York, NY, USA, 2008. ACM. ISBN 978-1-60558-011-1. doi: 10.1145/1357054.1357127. URL http://doi.acm.org/10.1145/1357054.1357127.

[117] C. Koch and S. Ullman. Shifts in selective visual attention: towards the underlying neural circuitry. *Human Neurbiology*, 4:219–227, 1985.

[118] Kathryn Koehler, Fei Guo, Sheng Zhang, and Miguel P. Eckstein. What do saliency models predict? *Journal of Vision*, 14(3):14, 2014. doi: 10.1167/14.3.14. URL +http://dx.doi.org/10.1167/14.3.14.

[119] Christof Körner. Eye movements reveal distinct search and reasoning processes in comprehension of complex graphs. *Applied Cognitive Psychology*, 25(6):893–905, 2011. ISSN 1099-0720. doi: 10.1002/acp.1766. URL http://dx.doi.org/10.1002/acp.1766.

[120] Stephen M Kosslyn. Understanding charts and graphs. *Applied Cognitive Psy-chology*, 3(3):185–225, 1989.

[121] A. Kovashka and J. Hahn. Towards automatic understanding of visual adver-tisements. In *CVPR Workshops*, 2018. URL http://people.cs.pitt.edu/~kovashka/ads_workshop/.

[122] Eileen Kowler. The role of visual and cognitive processes in the control of eye movement. *Reviews of oculomotor research*, 4:1–70, 1989.

[123] Kyle Krafka, Aditya Khosla, Petr Kellnhofer, Harini Kannan, Suchendra Bhandarkar, Wojciech Matusik, and Antonio Torralba. Eye tracking for everyone. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2176–2184, June 2016. doi: 10.1109/CVPR.2016.239.

[124] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in Neural Information Processing Systems*, pages 1097–1105, 2012.

[125] Srinivas S. Kruthiventi, Kumar Ayush, and R. Venkatesh Babu. Deepfix: A fully convolutional neural network for predicting human eye fixations. *CoRR*, abs/1510.02927, 2015. URL http://arxiv.org/abs/1510.02927.

[126] Ranjitha Kumar, Jerry O. Talton, Salman Ahmad, and Scott R. Klemmer. Bricolage: Example-based retargeting for web design. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, CHI '11, pages 2197–2206, New York, NY, USA, 2011. ACM. ISBN 978-1-4503-0228-9. doi: 10.1145/1978942. 1979262. URL http://doi.acm.org/10.1145/1978942.1979262.

[127] Ranjitha Kumar, Arvind Satyanarayan, Cesar Torres, Maxine Lim, Salman Ahmad, Scott R. Klemmer, and Jerry O. Talton. Webzeitgeist: Design mining the web. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, CHI '13, pages 3083–3092, New York, NY, USA, 2013. ACM. ISBN 978-1-4503-1899-0. doi: 10.1145/2470654.2466420. URL http://doi.acm.org/10.1145/2470654.2466420.

[128] M. Kümmerer, T. Wallis, and M. Bethge. How close are we to understanding image-based saliency? *arXiv preprint arXiv:1409.7686*, 2014.

[129] Matthias Kümmerer, Lucas Theis, and Matthias Bethge. Deep gaze I: boosting saliency prediction with feature maps trained on imagenet. *CoRR*, abs/1411.1045, 2014. URL http://arxiv.org/abs/1411.1045.

[130] Matthias Kümmerer, Thomas SA Wallis, and Matthias Bethge. Information-theoretic model comparison unifies saliency metrics. *Proceedings of the National Academy of Sciences*, 112(52):16054–16059, 2015.

[131] Matthias Kümmerer, Thomas S. A. Wallis, and Matthias Bethge. Saliency benchmarking: Separating models, maps and metrics. *CoRR*, abs/1704.08615, 2017. URL http://arxiv.org/abs/1704.08615.

[132] Dmitry Lagun and Eugene Agichtein. Viewser: Enabling large-scale remote user studies of web search examination and interaction. In *Proceedings of the 34th International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR '11, pages 365–374, New York, NY, USA, 2011.

ACM. ISBN 978-1-4503-0757-4. doi: 10.1145/2009916.2009967. URL http://doi.acm.org/10.1145/2009916.2009967.

[133] Guanbin Li and Yizhou Yu. Visual saliency based on multiscale deep features. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5455–5463, 2015.

[134] Huiyang Li and Nadine Moacdieh. Is chart junk useful? an extended examination of visual embellishment. In *Proceedings of the Human Factors and Ergonomics Society Annual Meeting*, volume 58, pages 1516–1520. SAGE Publications, 2014.

[135] Xi Li, Liming Zhao, Lina Wei, Ming-Hsuan Yang, Fei Wu, Yueting Zhuang, Haibin Ling, and Jingdong Wang. Deepsaliency: Multi-task deep neural network model for salient object detection. *IEEE Transactions on Image Processing*, 25(8):3919 – 3930, Aug 2016. ISSN 1057-7149. doi: 10.1109/TIP.2016.2579306.

[136] Daniel J. Liebling and Sören Preibusch. Privacy considerations for a pervasive eye tracking world. In *Proceedings of the 2014 ACM International Joint Conference on Pervasive and Ubiquitous Computing: Adjunct Publication*, UbiComp '14 Adjunct, pages 1169–1177, New York, NY, USA, 2014. ACM. ISBN 978-1-4503-3047-3. doi: 10.1145/2638728.2641688. URL http://doi.acm.org/10.1145/2638728.2641688.

[137] Sharon Lin and Pat Hanrahan. Modeling how people extract color themes from images. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, CHI '13, pages 3101–3110, New York, NY, USA, 2013. ACM. ISBN 978-1-4503-1899-0. doi: 10.1145/2470654.2466424. URL http://doi.acm.org/10.1145/2470654.2466424.

[138] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and Lawrence Zitnick. Microsoft coco: Common objects in context. In David Fleet, Tomas Pajdla, Bernt Schiele, and Tinne Tuytelaars, editors, *Computer Vision – ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V*, pages 740–755. Springer International Publishing, Cham, 2014. ISBN 978-3-319-10602-1. doi: 10.1007/978-3-319-10602-1_48. URL http://dx.doi.org/10.1007/978-3-319-10602-1_48.

[139] Nian Liu, Junwei Han, Dingwen Zhang, Shifeng Wen, and Tianming Liu. Predicting eye fixations using convolutional neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 362–370, 2015.

[140] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg. SSD: Single Shot MultiBox Detector. In *ECCV*, 2016.

[141] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 39(4):640–651, April 2017. ISSN 0162-8828. doi: 10.1109/TPAMI.2016.2572683.

[142] Christof Lutteroth, Moiz Penkar, and Gerald Weber. Gaze vs. mouse: A fast and accurate gaze-only click alternative. In *Proceedings of the 28th Annual ACM Symposium on User Interface Software &#38; Technology*, UIST '15, pages 385–394, New York, NY, USA, 2015. ACM. ISBN 978-1-4503-3779-3. doi: 10.1145/2807442.2807461. URL http://doi.acm.org/10.1145/2807442.2807461.

[143] Päivi Majaranta and Andreas Bulling. *Eye Tracking and Eye-Based Human–Computer Interaction*, pages 39–65. Springer London, London, 2014. ISBN 978-1-4471-6392-3. doi: 10.1007/978-1-4471-6392-3_3. URL http://dx.doi.org/10.1007/978-1-4471-6392-3_3.

[144] Simone Marinai, Beatrice Miotti, and Giovanni Soda. Digital libraries and document image retrieval techniques: A survey. In *Learning Structure and Schemas from Documents*, pages 181–204. Springer, 2011.

[145] Kim Marriott, Helen Purchase, Michael Wybrow, and Cagatay Goncu. Memorability of visual features in network diagrams. *IEEE TVCG*, 18(12):2477–2485, 2012.

[146] George W. McConkie and Keith Rayner. The span of the effective stimulus during a fixation in reading. *Perception & Psychophysics*, 17(6):578–586, 1975. ISSN 1532-5962. doi: 10.3758/BF03203972. URL http://dx.doi.org/10.3758/BF03203972.

[147] O. Le Meur and T. Baccino. Methods for comparing scanpaths and saliency maps: strengths and weaknesses. *Behavioral Research Methods*, 45(1):251–266, 2013.

[148] O. Le Meur, P. Le Callet, and D. Barba. Predicting visual fixations on video based on low-level visual features. *Vision Research*, 47(19):2483–2498, 2007.

[149] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space. *CoRR*, abs/1301.3781, 2013. URL http://arxiv.org/abs/1301.3781.

[150] George A Miller. WordNet: a lexical database for English. *Communications of the ACM*, 38(11):39–41, 1995.

[151] Jakob Nielsen and Kara Pernice. *Eyetracking Web Usability*. New Riders Publishing, Thousand Oaks, CA, USA, 1st edition, 2009. ISBN 0321498364, 9780321498366.

[152] David Noton and Lawrence Stark. Scanpaths in saccadic eye movements while viewing and recognizing patterns. *Vision Research*, 11(9):929 – IN8, 1971. ISSN 0042-6989. doi: https://doi.org/10.1016/0042-6989(71)90213-6. URL http://www.sciencedirect.com/science/article/pii/0042698971902136.

[153] Peter O'Donovan, Aseem Agarwala, and Aaron Hertzmann. Color compatibility from large datasets. *ACM TOG*, 30(43), 2011.

[154] Peter O'Donovan, Aseem Agarwala, and Aaron Hertzmann. Learning layouts for single-page graphic designs. *IEEE Transactions on Visualization and Computer Graphics*, 20(8):1200–1213, Aug 2014. ISSN 1077-2626. doi: 10.1109/TVCG.2014.48.

[155] Peter O'Donovan, Aseem Agarwala, and Aaron Hertzmann. Designscape: Design with interactive layout suggestions. In *Proceedings of the 33rd Annual ACM Conference on Human Factors in Computing Systems*, CHI '15, pages 1221–1224, New York, NY, USA, 2015. ACM. ISBN 978-1-4503-3145-6. doi: 10.1145/2702123.2702149. URL http://doi.acm.org/10.1145/2702123.2702149.

[156] Bing Pan, Helene A. Hembrooke, Geri K. Gay, Laura A. Granka, Matthew K. Feusner, and Jill K. Newman. The determinants of web page viewing behavior: An eye-tracking study. In *Proceedings of the 2004 Symposium on Eye Tracking Research & Applications*, ETRA '04, pages 147–154, New York, NY, USA, 2004. ACM. ISBN 1-58113-825-3. doi: 10.1145/968363.968391. URL http://doi.acm.org/10.1145/968363.968391.

[157] Junting Pan, Elisa Sayrol, Xavier Giro-i Nieto, Kevin McGuinness, and Noel E O'Connor. Shallow and deep convolutional networks for saliency prediction. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 598–606, 2016.

[158] Anshul Vikram Pandey, Katharina Rall, Margaret L Satterthwaite, Oded Nov, and Enrico Bertini. How deceptive are deceptive visualizations?: An empirical analysis of common distortion techniques. In *CHI '15*, pages 15–03, 2015.

[159] Xufang Pang, Ying Cao, Rynson W. H. Lau, and Antoni B. Chan. Directing user attention via visual flow on web designs. *ACM Trans. Graph.*, 35(6):240:1–240:11, November 2016. ISSN 0730-0301. doi: 10.1145/2980179.2982422. URL http://doi.acm.org/10.1145/2980179.2982422.

[160] Alexandra Papoutsaki, Patsorn Sangkloy, James Laskey, Nediyana Daskalova, Jeff Huang, and James Hays. Webgazer: Scalable webcam eye tracking using user interactions. In *Proceedings of the 25th International Joint Conference on Artificial Intelligence (IJCAI)*, pages 3839–3845. AAAI, 2016.

[161] D. Parkhurst, K. Law, and E. Niebur. Modeling the role of salience in the allocation of overt visual attention. *Vision Research*, 42(1):107 – 123, 2002. ISSN 0042-6989. doi: DOI:10.1016/S0042-6989(01)00250-4. URL http://www.sciencedirect.com/science/article/B6T0W-44XM16S-8/2/0cf0f73ccc1140e2bf9615208c10f1c0.

[162] O. Pele and M. Werman. A linear time histogram metric for improved sift matching. In *ECCV*, 2008.

[163] R. J. Peters, A. Iyer, L. Itti, and C. Koch. Components of bottom-up gaze allocation in natural images. *Vision Research*, 45(18):2397 – 2416, 2005. ISSN 0042-6989. doi: DOI:10.1016/j.visres.2005.03.019. URL http://www.sciencedirect.com/science/article/B6T0W-4G9GN35-1/2/a0353bdde8c613f899bdb78568a6f763.

[164] P. O. Pinheiro, T.-Y. Lin, R. Collobert, and P. Dollàr. Learning to Refine Object Segments. In *ECCV*, 2016.

[165] Pedro O. Pinheiro, Ronan Collobert, and Piotr Dollr. Learning to segment object candidates. In *NIPS*, 2015.

[166] Steven Pinker. A theory of graph comprehension. *Artificial intelligence and the future of testing*, pages 73–126, 1990.

[167] Jorge Poco and Jeffrey Heer. Reverse-engineering visualizations: Recovering visual encodings from chart images. *Computer Graphics Forum (Proc. EuroVis)*, 2017. URL http://idl.cs.washington.edu/papers/reverse-engineering-vis.

[168] Mathias Pohl, Markus Schmitt, and Stephan Diehl. Comparing the readability of graph layouts using eyetracking and task-oriented analysis. In *Proceedings of the Fifth Eurographics Conference on Computational Aesthetics in Graphics, Visualization and Imaging*, Computational Aesthetics'09, pages 49–56, Aire-la-Ville, Switzerland, Switzerland, 2009. Eurographics Association. ISBN 978-3-905674-17-0. doi: 10.2312/COMPAESTH/COMPAESTH09/049-056. URL http://dx.doi.org/10.2312/COMPAESTH/COMPAESTH09/049-056.

[169] Alex Poole and Linden J Ball. Eye tracking in HCI and usability research. *Encyclopedia of human computer interaction*, 1:211–219, 2006.

[170] V Shiv Naga Prasad, Behjat Siddiquie, Jennifer Golbeck, and Larry S Davis. Classifying computer generated charts. In *CBMI*, pages 85–92. IEEE, 2007.

[171] Michael Raschke, Tanja Blascheck, Marianne Richter, Tanja Agapkin, and Thomas Ertl. Visual analysis of perceptual and cognitive processes. In *IJCV*, 2014.

[172] Keith Rayner. Eye movements in reading and information processing: 20 years of research. *Psychological bulletin*, 124(3):372, 1998.

[173] Keith Rayner. The gaze-contingent moving window in reading: Development and review. *Visual Cognition*, 22(3-4):242–258, 2014.

[174] Keith Rayner, Caren M Rotello, Andrew J Stewart, Jessica Keir, and Susan A Duffy. Integrating text and pictorial information: eye movements when looking at print advertisements. *Journal of Experimental Psychology: Applied*, 7(3):219, 2001.

[175] Ali Sharif Razavian, Hossein Azizpour, Josephine Sullivan, and Stefan Carlsson. Cnn features off-the-shelf: An astounding baseline for recognition. In *Proceedings of the 2014 IEEE Conference on Computer Vision and Pattern Recognition Workshops*, CVPRW '14, pages 512–519, Washington, DC, USA, 2014. IEEE Computer Society. ISBN 978-1-4799-4308-1. doi: 10.1109/CVPRW.2014.131. URL http://dx.doi.org/10.1109/CVPRW.2014.131.

[176] Adria Recasens, Aditya Khosla, Carl Vondrick, and Antonio Torralba. Where are they looking? In *Advances in Neural Information Processing Systems*, pages 199–207, 2015.

[177] J. Redmon and A. Farhadi. YOLO9000: Better, Faster, Stronger. In *CVPR*, 2017.

[178] Eyal M Reingold, Lester C Loschky, George W McConkie, and David M Stampe. Gaze-contingent multiresolutional displays: An integrative review. *Human Factors: The Journal of the Human Factors and Ergonomics Society*, 45(2):307–328, 2003.

[179] S. Ren, K. He, R. Girshick, and J. Sun. Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks. In *ICCV*, 2015.

[180] Ronald A Rensink. *The management of visual attention in graphic displays*. Cambridge University Press, Cambridge, England, 2011.

[181] N. Riche, M. Duvinage, M. Mancas, B. Gosselin, and T. Dutoit. Saliency and human fixations: State-of-the-art and study of comparison metrics. In *ICCV*, 2013.

[182] Stephan R. Richter, Vibhav Vineet, Stefan Roth, and Vladlen Koltun. Playing for data: Ground truth from computer games. In *ECCV*, 2016.

[183] Daniel Ritchie, Ankita Arvind Kejriwal, and Scott R Klemmer. d. tour: Style-based exploration of design example galleries. In *UIST*, pages 165–174. ACM, 2011.

[184] Kerry Rodden, Xin Fu, Anne Aula, and Ian Spiro. Eye-mouse coordination patterns on web search results pages. In *CHI '08 Extended Abstracts on Human Factors in Computing Systems*, CHI EA '08, pages 2997–3002, New York, NY, USA, 2008. ACM. ISBN 978-1-60558-012-8. doi: 10.1145/1358628.1358797. URL http://doi.acm.org/10.1145/1358628.1358797.

[185] Ruth Rosenholtz, Amal Dorai, and Rosalind Freeman. Do predictions of visual perception aid design? *ACM Trans. Appl. Percept.*, 8(2):12:1–12:20, February 2011. ISSN 1544-3558. doi: 10.1145/1870076.1870080. URL http://doi.acm.org/10.1145/1870076.1870080.

[186] Michael Rubinstein, Diego Gutierrez, Olga Sorkine, and Ariel Shamir. A comparative study of image retargeting. *ACM Trans. Graph.*, 29(6):160:1–160:10, December 2010. ISSN 0730-0301. doi: 10.1145/1882261.1866186. URL http://doi.acm.org/10.1145/1882261.1866186.

[187] Y. Rubner, C. Tomasi, and L. J. Guibas. The earth movers distance as a metric for image retrieval. *IJCV*, 40, 2000.

[188] Babak Saleh, Mira Dontcheva, Aaron Hertzmann, and Zhicheng Liu. Learning style similarity for searching infographics. In *Proceedings of the 41st graphics interface conference*, pages 59–64. Canadian Information Processing Society, 2015.

[189] Patsorn Sangkloy, Nathan Burnell, Cusuh Ham, and James Hays. The sketchy database: learning to retrieve badly drawn bunnies. *ACM Transactions on Graphics (TOG)*, 35(4):119, 2016.

[190] Guy Satat, Matthew Tancik, Otkrist Gupta, Barmak Heshmat, and Ramesh Raskar. Object classification through scattering media with deep learning on time resolved measurement. *Optics Express*, 25(15):17466–17479, 2017.

[191] Manolis Savva, Nicholas Kong, Arti Chhajta, Li Fei-Fei, Maneesh Agrawala, and Jeffrey Heer. Revision: Automated classification, analysis and redesign of chart images. In *Proceedings of the 24th Annual ACM Symposium on User Interface Software and Technology*, UIST '11, pages 393–402, New York, NY, USA, 2011. ACM. ISBN 978-1-4503-0716-1. doi: 10.1145/2047196.2047247. URL http://doi.acm.org/10.1145/2047196.2047247.

[192] Michael Schulte-Mecklenbeck, Ryan O. Murphy, and Florian Hutzler. Flashlight - recording information acquisition online. *Comput. Hum. Behav.*, 27(5):1771–1782, September 2011. ISSN 0747-5632. doi: 10.1016/j.chb.2011.03.004. URL http://dx.doi.org/10.1016/j.chb.2011.03.004.

[193] Min Joon Seo, Hannaneh Hajishirzi, Ali Farhadi, and Oren Etzioni. Diagram understanding in geometry questions. In *AAAI*, pages 2831–2838, 2014.

[194] Chengyao Shen and Qi Zhao. *Webpage Saliency*, pages 33–46. Springer International Publishing (ECCV), Cham, 2014. ISBN 978-3-319-10584-0. doi: 10.1007/978-3-319-10584-0_3. URL http://dx.doi.org/10.1007/978-3-319-10584-0_3.

[195] Chengyao Shen, Xun Huang, and Qi Zhao. Predicting eye fixations on webpage with an ensemble of early features and high-level representations from deep network. *IEEE Transactions on Multimedia*, 17(11):2084–2093, Nov 2015. ISSN 1520-9210. doi: 10.1109/TMM.2015.2483370.

[196] Noah Siegel, Zachary Horvitz, Roie Levin, Santosh Divvala, and Ali Farhadi. Figureseer: Parsing result-figures in research papers. In *European Conference on Computer Vision*, pages 664–680. Springer, 2016.

[197] H. Siirtola, T. Laivo, T. Heimonen, and K.-J. Raiha. Visual perception of parallel coordinate visualizations. In *Intern. Conf. on Information Visualisation*, pages 3–9, July 2009. doi: 10.1109/IV.2009.25.

[198] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *CoRR*, abs/1409.1556, 2014. URL http://arxiv.org/abs/1409.1556.

[199] Drew Skau, Lane Harrison, and Robert Kosara. An evaluation of the impact of visual embellishments in bar charts. In *Proceedings of EuroVis*, volume 2015, 2015.

[200] Ray Smith. An overview of the tesseract ocr engine. In *ICDAR*, volume 2, pages 629–633. IEEE, 2007.

[201] Hyun Soo Park and Jianbo Shi. Social saliency prediction. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4777–4785, 2015.

[202] Jeremiah D Still and Christopher M Masciocchi. A saliency model predicts fixations in web interfaces. In *5 th International Workshop on Model Driven Development of Advanced User Interfaces (MDDAUI 2010)*, page 25. Citeseer, 2010.

[203] Hendrik Strobelt, Daniela Oelke, Christian Rohrdantz, Andreas Stoffel, Daniel A Keim, and Oliver Deussen. Document cards: A top trumps visualization for documents. *IEEE Transactions on Visualization and Computer Graphics*, 15(6):1145–1152, 2009.

[204] M. J. Swain and D. H. Ballard. Color indexing. *IJCV*, 7(1):11–32, 1991.

[205] Peter Tarasewich, Marc Pomplun, Stephanie Fillion, and Daniel Broberg. The enhanced restricted focus viewer. *International Journal of HumanComputer Interaction*, 19(1):35–54, 2005. doi: 10.1207/s15327590ijhc1901\_4.

[206] Benjamin W. Tatler. The central fixation bias in scene viewing: Selecting an optimal viewing position independently of motor biases and image feature distributions. *Journal of Vision*, 7(14):4, 2007. doi: 10.1167/7.14.4. URL +http://dx.doi.org/10.1167/7.14.4.

[207] Benjamin W. Tatler, Roland J. Baddeley, and Iain D. Gilchrist. Visual correlates of fixation selection: effects of scale and time. *Vision Research*, 45(5):643 – 659, 2005. ISSN 0042-6989. doi: https://doi.org/10.1016/j.visres.2004.09.017. URL http://www.sciencedirect.com/science/article/pii/S0042698904004626.

[208] Benjamin W. Tatler, Mary M. Hayhoe, Michael F. Land, and Dana H. Ballard. Eye guidance in natural vision: Reinterpreting salience. *Journal of Vision*, 11(5): 5, 2011. doi: 10.1167/11.5.5. URL +http://dx.doi.org/10.1167/11.5.5.

[209] Hamed R Tavakoli, Fawad Ahmed, Ali Borji, and Jorma Laaksonen. Saliency revisited: Analysis of mouse movements versus fixations. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017.

[210] Jaime Teevan, Edward Cutrell, Danyel Fisher, Steven M. Drucker, Gonzalo Ramos, Paul André, and Chang Hu. Visual snippets: Summarizing web pages for search and revisitation. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, CHI '09, pages 2023–2032, New York, NY, USA, 2009. ACM. ISBN 978-1-60558-246-7. doi: 10.1145/1518701.1519008. URL http://doi.acm.org/10.1145/1518701.1519008.

[211] Tobii. Tobii eye tracking: An introduction to eye tracking and tobii eye trackers. White paper, Tobii Technology AB, 2010.

[212] A. M. Treisman and G. Gelade. A feature-integration theory of attention. *Cognitive Psychology*, 12(1):97 – 136, 1980. ISSN 0010-0285. doi: DOI:10. 1016/0010-0285(80)90005-5. URL http://www.sciencedirect.com/science/article/B6WCR-4D6RJM2-46/2/b34aa05384b2a7702189c22840489174.

[213] Satoshi Tsutsui and David Crandall. A data driven approach for compound figure separation using convolutional neural networks. *ICDAR*, 2017.

[214] Andrew Vande Moere, Martin Tomitsch, Christoph Wimmer, Boesch Christoph, and Thomas Grechenig. Evaluating the effect of style in information visualization. *IEEE TVCG*, 18(12):2739–2748, 2012.

[215] Eleonora Vig, Michael Dorr, and David Cox. Large-scale optimization of hierarchical features for saliency prediction in natural images. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2798–2805, 2014.

[216] Luis von Ahn and Laura Dabbish. Labeling images with a computer game. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, CHI '04, pages 319–326, New York, NY, USA, 2004. ACM. ISBN 1-58113-702-8. doi: 10.1145/985692.985733. URL http://doi.acm.org/10.1145/985692.985733.

[217] Lijun Wang, Huchuan Lu, Xiang Ruan, and Ming-Hsuan Yang. Deep networks for saliency detection via local estimation and global search. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3183–3192, 2015.

[218] Michael J Wilber, Chen Fang, Hailin Jin, Aaron Hertzmann, John Collomosse, and Serge Belongie. BAM! The Behance Artistic Media Dataset for Recognition Beyond Photography. *ICCV*, 2017.

[219] N. Wilming, T. Betz, T. C. Kietzmann, and P. König. Measures and limits of models of fixation selection. *PLoS ONE*, 6, 2011.

[220] Allison Woodruff, Andrew Faulring, Ruth Rosenholtz, Julie Morrsion, and Peter Pirolli. Using thumbnails to search the web. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, CHI '01, pages 198–205, New York, NY, USA, 2001. ACM. ISBN 1-58113-327-8. doi: 10.1145/365024.365098. URL http://doi.acm.org/10.1145/365024.365098.

[221] Peng Wu, Sandra Carberry, Stephanie Elzer, and Daniel Chester. Recognizing the intended message of line graphs. In *Diagrammatic Representation and Inference*, pages 220–234. Springer, 2010.

[222] Jianxiong Xiao, James Hays, Krista A Ehinger, Aude Oliva, and Antonio Torralba. Sun database: Large-scale scene recognition from abbey to zoo. In *IEEE conference on Computer vision and pattern recognition (CVPR)*, pages 3485–3492, 2010.

[223] J. Xu, M. Jiang, S. Wang, M. S. Kankanhalli, and Q. Zhao. Predicting human gaze beyond pixels. *Journal of Vision*, 14(1):1–20, 2014.

[224] Pingmei Xu, Krista A. Ehinger, Yinda Zhang, Adam Finkelstein, Sanjeev R. Kulkarni, and Jianxiong Xiao. Turkergaze: Crowdsourcing saliency with webcam based eye tracking. *CoRR*, abs/1504.06755, 2015. URL http://arxiv.org/abs/1504.06755.

[225] Pingmei Xu, Yusuke Sugano, and Andreas Bulling. Spatio-temporal modeling and prediction of visual attention in graphical user interfaces. In *Proceedings of the 2016 CHI Conference on Human Factors in Computing Systems*, CHI '16, pages 3299–3310, New York, NY, USA, 2016. ACM. ISBN 978-1-4503-3362-7. doi: 10.1145/2858036.2858479. URL http://doi.acm.org/10.1145/2858036.2858479.

[226] Bangpeng Yao, Xiaoye Jiang, Aditya Khosla, Andy Lai Lin, Leonidas Guibas, and Li Fei-Fei. Human action recognition by learning bases of action attributes and parts. In *IEEE International Conference on Computer Vision (ICCV)*, pages 1331–1338, 2011.

[227] J. Zhang and S. Sclaroff. Saliency detection: a boolean map approach. In *IEEE International Conference on Computer Vision*, 2013.

[228] L. Zhang, M. H. Tong, T. K. Marks, H. Shan, and G. W. Cottrell. Sun: A bayesian framework for saliency using natural statistics. *Journal of Vision*, 8(7), 2008. doi: 10.1167/8.7.32.

[229] Xucong Zhang, Yusuke Sugano, Mario Fritz, and Andreas Bulling. Appearance-based gaze estimation in the wild. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4511–4520, 2015.

[230] Qi Zhao and Christof Koch. Learning a saliency map using fixated locations in natural scenes. *Journal of Vision*, 11(3):9–9, 2011.

[231] Rui Zhao, Wanli Ouyang, Hongsheng Li, and Xiaogang Wang. Saliency detection by multi-context deep learning. In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1265–1274, June 2015. doi: 10.1109/ CVPR.2015.7298731.

[232] C Lawrence Zitnick and Devi Parikh. Bringing semantics into focus using visual abstraction. In *CVPR*, 2013.

[233] C Lawrence Zitnick, Ramakrishna Vedantam, and Devi Parikh. Adopting abstract images for semantic scene understanding. *IEEE TPAMI*, 38(4):627–638, 2016.