# Influence of Vertical Result in Web Search Examination

Zeyang Liu[†], Yiqun Liu[†][*], Ke Zhou[‡], Min Zhang[†], Shaoping Ma[†]

[†]Tsinghua National Laboratory for Information Science and Technology, Department of Computer Science & Technology, Tsinghua University, Beijing, China

[‡]Yahoo Labs, London, U.K.

yiqunliu@tsinghua.edu.cn

## ABSTRACT

Research in how users examine results on search engine result pages (SERPs) helps improve result ranking, advertisement placement, performance evaluation and search UI design. Although examination behavior on organic search results (also known as "ten blue links") has been well studied in existing works, there lacks a thorough investigation on how users examine SERPs with verticals. Considering the fact that a large fraction of SERPs are served with one or more verticals in the practical Web search scenario, it is of vital importance to understand the influence of vertical results on search examination behaviors. In this paper, we focus on five popular vertical types and try to study their influences on users' examination processes in both cases when they are relevant or irrelevant to the search queries. With examination behavior data collected with an eye-tracking device, we show the existence of vertical-aware user behavior effects including vertical *attraction* effect, examination *cut-off* effect in the presence of a relevant vertical, and examination *spill-over* effect in the presence of an irrelevant vertical. Furthermore, we are also among the first to systematically investigate the internal examination behavior within the vertical results. We believe that this work will promote our understanding of user interactions with federated search engines and bring benefit to the construction of search performance evaluations.

## Categories and Subject Descriptors

H.3.3 [**Information Storage and Retrieval**]: Information Search and Retrieval

## Keywords

Federated search; User behavior analysis; Eye tracking

## 1. INTRODUCTION

With the rapid evolvement of Web search engines, traditional search engine result pages (SERPs), which consist of ten web documents (known as ten "blue links" or "organic results"), might become not so effective when users issue their queries to search engines. To advance the efficiency of information search behaviors, it is popular for commercial search portals such as Google, Yahoo! to blend some specific vertical results, whose contents are assembled from other data sources, into the ranked list of standard results. Since users could instantly obtain the information that they expected from relevant verticals, search behavior might be inevitably affected by these "special results". A good example of this is attractiveness bias [5, 9, 28], which shows that users may adopt their visual attention to verticals more quickly, ignoring the other documents on SERPs. It is, therefore, important to gain a good understanding of the influences of vertical results on search behaviors in designing federated search system.

Because of the special design of vertical documents, the embedment of verticals has revolutionized the way users find information on SERPs. Many prior studies have revealed this phenomenon. Chen et al [9] are among the first to propose a federated click model and demonstrate the difference in users' click behaviors in federated search. Wang et al. [28] further found the existence of presentation bias that may lead to different examination behavior on both verticals and other surrounding results, with the help of eye-tracking equipments. Lagun et al. [20] and Arguello et al. [1] indicated that the relevance of verticals could affect search behaviors, including gaze activity and cursor movements.

Although these previous works have revealed the difference in user behavior between federated and traditional search, it is still not sufficient to draw a complete picture of the process of user examination in practical web search scenario. Firstly, compared to the homogeneous presentation style of the organic results, verticals have a wide variety of presentation styles, which might lead to different examination behaviors on the SERPs. Most prior studies, such as [5, 9, 28, 29], commonly divided the vertical type into two categories: textual and multimedia. And the effects of these two types were measured and analyzed independently. However, with the increasing diversity of verticals, the content of vertical documents does not only contain one single type of items. Rather, the vertical blocks that contain both multimedia and textual snippets (e.g., news) exist widely in commercial search engines. *The examination behavior might vary on this "mixed type verticals".* For example, Figure 1 shows an example of the attention heatmaps for news and image verticals. Interestingly, while users are more likely to be attracted by multimedia verticals [9, 28], our findings indicate that news vertical, which also contained image element in their snippets, do not attract more visual attention on the image (i.e. the users examine carefully on

the textual snippets without much attention on the image snippet). Secondly, previous works [1, 18, 19] pointed out that vertical quality leads to different user interaction patterns and satisfactions. However, these studies only focus on the page level effect, such as the spill-over effect [1], which means that users might prefer to interact with the web results when the vertical corresponds to the users' intention. To some extent, *this limits the in-depth analysis of user decision process at result level within verticals of different quality.* Essentially, as we demonstrate in our experiments in Section 4.1, the spill-over effect may contribute to different examination distributions on organic results, especially those close to the verticals. In addition, very little is known about *the internal examination of different items presented in the vertical block* from the previous studies. A better understanding of the internal examination will help us take a step further to comprehend the cognitive mechanisms in federated search and design better vertical user interfaces.



**Figure 1: Users' attention heatmaps for news and image vertical on two SERPs (Left: news; Right: image.)**

The contribution of the paper is two fold:

- We study the influence of vertical type, embedding position and vertical relevance on users' search examination process;
- We investigate the examination process among different sub-components within the vertical block.

The rest of the paper is organized as follows. Section 2 presents an overview of the related work. Then we thoroughly describe the experimental design for collecting user behavior data. In section 4, we provide the analysis of user examination behavior on SERPs. A summary of this paper is presented in Section 5.

## 2. RELATED WORK

### 2.1 Federated Search

How to federate search results from vertical sources into organic web search results has been extensively studied in recent years [1, 9, 11, 28, 3, 2, 21]. Most of the current researches focus on how to select the most relevant verticals [3, 2, 21], how to appropriately embed vertical results into organic web results [3, 2] and how to measure the effectiveness of the federated search pages [9, 28, 29]. Federation on the web can be also referred as aggregated search [1, 9, 28, 29].

With the introduction of verticals, the user behavior has become more and more complex. Recently, several work aim to better understand users' search behavior by either conducting user studies or performing large-scale log analysis. Sushmita et al. [27] conducted a user study and found that users click more on video results than the image and news results. Zhou et al. [30]

performed a crowd-sourcing study with explicit user assessments and found similar visual saliency biases. Diaz et al. [11] mined users' mouse movement interactions and found that different result appearances might lead to different bias strengths. Arguello et al. [1] focused on studying aggregated search coherence, which means the extent to which results from different verticals focus on similar senses of an ambiguous or underspecified query. They found that users are more likely to interact with the web results when the vertical results are consistent with the user's intended query sense. By exploiting a commercial search engine log, Chen et al. [9] found that users are more likely to end his/her search sessions immediately after clicking the vertical results.

Most of the work mentioned above focus on implicit signals such as clicks and mouse movements, or explicit user assessments while we move one step forward, towards further understanding of user behavior (especially users' attention) through eye-tracking devices. We revisit some of the biases such as visual saliency, vertical relevance and vertical ranked position, and also obtain several new findings.

### 2.2 Eye-tracking Studies in Web Search

Eye-tracking device has been widely utilized in understanding user behavior on both sponsored search [8, 24, 12] and organic search results [14, 18, 19]. Since the eye tracker can record users' real-time eye movement information on SERPs, it helps researchers better understand how users examine results. Granka et al. [14], and Joachims et al. [18, 19] are among the first works that start this line of research on organic SERPs and they all found that there exists position bias during users' examination processes. [10] revealed the connection between the length of contextual snippet and user performance by using eye tracking techniques. Recently, Liu et al. [22] indicated that the examination process of search results might have two steps (skim and read) and proposed a two-stage examination model based on the eye-tracking analysis. Eye-tracking experiments also help analyze the users' evaluation on search results. For example, [4] showed that the personal style (economic or exhaustive) and the result relevance can affect search result evaluation. [7] predict the relevance of documents by using the gaze data. Besides, eye-tracking devices are used to investigate the relationship between eye movements and mouse movements, such as [25, 26, 15, 17]

Researchers also conducted eye-tracking experiments in the context of federated search. Wang et al. [28] found that different verticals might create examination biases on the eye movement behavior for both vertical and other results on SERPs. Lagun et al. [20] performed a similar study on mobile devices and found that both vertical relevance and positions have impacts on users' attentions. [13] investigated the influence of social annotations, which present social signals or information, on users' performance in web search. Navalpakkam et al. [23] found that the flow of user attention on non-linear page layouts (e.g., with the existence of knowledge graph components) is different from the widely believed top-down linear examination order of the search results.

Compared to previous work, our work performs a more extensive eye-tracking study that examines a variety of vertical biases and their impacts in organic results on the SERP. We also provide insights for verticals (with five different presentation styles). Another contribution is that we study in-depth users' examination behaviors on different items shown within vertical blocks, which has not been studied before. The insights provided can be helpful in assisting vertical-aware click models [9, 28] by better estimating item-level examination probability.

# 3. METHODS AND MATERIALS

## 3.1 Experimental Design

To gain good insights into users' search strategies and processes on SERPs with verticals, we designed an eye-tracking experiment in which participants were asked to complete 30 search tasks. Considering that different layouts and presentation styles of the SERPs may affect users' behavior, we implemented a controlled search engine system to collect users' interaction data. With the system, we present different SERPs to users and track their interaction, such as the gaze duration or movement patterns on the different results of the SERPs.

The SERPs we present to the users vary in three aspects that might have effects in users' examination behavior (i.e. the three independent variables we manipulated in the study):

- Vertical type (see Figure 4(b) for examples): textual, encyclopedia, image-only, application-download, news or none (i.e. organic results). Due to the existence of presentation bias in federated search [9, 28], vertical type plays an important role in affecting users' search behavior. In our experimentation, we selected five different types, which are among the most popular types in commercial search engines, and embedded them into the organic result list respectively.
- Vertical position. Since position bias affects user attention in federated search engine [9, 28], this factor was also taken into account in our experimentation. Vertical results are randomly placed at position 1, 3 and 5 of result lists, respectively.
- Vertical quality. Users' interaction on SERPs may also be influenced by the relevance of results, not only in traditional "ten blue links" but also in sponsored [8] and federated searches [1]. Therefore, both "on topic" (relevant) and "off-topic" (irrelevant) results are also included in result lists.

The *dependent variables* we aim to track in terms of users' examination are:

- Vertical block itself. This again is to study the presentation bias of vertical results under different conditions as an extension of previous work [14, 28].
- Organic results around the vertical. This is to study the users' interaction of the organic results from the perspective of eye fixation under different conditions, which is different from previous works that focus only on implicit interaction signals [1, 11] (e.g. clicks, mouse movements).
- Items within vertical blocks. This is to investigate the examination patterns and the interaction of different elements in the vertical block, which is novel and not studied in previous work.

## 3.2 Experimental Protocol

Our experimental study is processed in the following steps. Firstly, to ensure that each participant was familiar with the operation of our experimental system and the experiment procedure, they were asked to finish two warm-up tasks. During this step, participants only perceived that their interactions with the search system, including eye activities and mouse movements, will be recorded and were unaware of the real purpose of our study. Then, participants were asked to go through calibration processes before they start the main search tasks so as to collect accurate and reliable eye movement data. Next, every participant performed the same set of 30 search tasks, followed by an exit questionnaire. To eliminate the influence of the different levels of task difficulty [17], we verified that the corresponding results list without verticals, namely organic results list, contained the answers to these tasks, and each task was moderately easy for most participants. At the end of the experiment, participants were required to give some feedback about their search experience and compensated with approximately US$10.

Before each task, participants were directed to a task description page and given an initial query for this task. In order to eliminate the possible ambiguities of the queries, the corresponding description of each query was also displayed on this page. It is worth noting that all of these 30 queries were selected from real-world commercial search logs so that they contain the practical users' search intention. Once participants had read the task description and confirmed what they should look for, they could click the search button and begin the search processes on the SERPs. To make sure that all participants see the same SERPs while performing a certaining search tasks, we crawled the SERPs from search engines and stored them in our web servers. This allowed us to have a consistent initial SERP for each query and to strictly control the experimental variables, which may lead to differences in users' behavior. After these initial SERPs were displayed, there were no restrictions on users' behavior on the search result pages. In other words, participants could be free to behave on the SERP as usual in all ways, such as click links or scroll the screen. The aim of this design is to create a realistic scenario for federated search in a laboratory environment. Furthermore, during the search process, participants' eye movements were recorded by the eye-tracker and their cursor activities, such as click, hover and scroll, were also simultaneously logged by the embedded JavaScript code. Once users are satisfied or too frustrated to continue, they could finish browsing the SERP and click the "finished" button to move on to the next search task. It took about 40 to 60 minutes for each participant to complete the whole experimentation.

### 3.2.1 SERP Generation

In this paper, we focus on the influence of five major types of vertical results, which were widely applied in current search engines. We modelled the layout of the SERPs based on a commercial search engine, which ensure the realistic and natural appearance of result lists. A federated SERP in our experimentation mainly consist of one specific vertical and nine organic links. All of these results, including vertical and organic results, were crawled from the same search engine, and the presentation order of these organic results remained the same. When a participant searched a query, the vertical result (if any) was randomly selected from a pool of verticals of different quality (relevant or irrelevant) and integrated into the SERP at Position 1 (the top of the first viewport), Position 3 (the middle of the first viewport) or Position 5 (the bottom of the first viewport). It is worth noting that all the organic result lists contained the answers (see Section 3.1) to the tasks so that the difficulty of the federated SERPs was similar regardless of the relevance of the vertical. Following the pre-procedure steps above, we implemented 180 initially fixed SERPs with verticals and 30 organic-only SERPs in total.

### 3.2.2 Apparatus

We deployed a Tobii X2-30 eye tracker to capture participants' eye movements and deploy the search system on a 17" LCD monitor whose resolution is 1360*768. The Internet Explorer 11 browser was used to display the pages of search system, including the description pages and search result pages. As the vertical may

occupy more space of the SERP, the page fold was commonly between the results at the positions 5 and 6. For identifying user examination behavior, we detect fixations using built-in algorithms from Tobii Studio. In these algorithms, the gaze whose duration was above 60ms around a specific location would be collected and treated as a fixation.

### 3.2.3 Participants

Altogether 35 participants (aged 18 to 25, mean = 18.8) with informed consent were recruited in our experimentation. Because subjects were expected to have a variety of backgrounds, all of these participants were selected from a wide range of majors (e.g., biology, economics, engineering etc.) of a university. Because of the calibration problems with the eye tracker, not all of their eye movement data were available. As a result, data from 32 of these participants was finally taken into account. Particularly, most of the participants had normal vision and were able to browse the on-screen web pages without wearing glasses. Additionally, all the participants self-reported that they were familiar with the operation of search engines and were confident that the SERPs that we provided come from the original search engine without modifications.

## 3.3 Search Tasks

Each of participants was instructed to perform 30 search tasks using our search system. As described in Section 3.1, to generate appropriate initial queries, we sample a set of "medium frequency" queries from raw search logs, which come from a major commercial search engine. And then, the original SERPs related to these queries were examined one by one, and the SERPs that contained less than three specific verticals could be selected into the initial SERP set. This setup was chosen to reduce the influence of excessive modification of SERPs and be consistent with the realistic scenario as far as possible. Finally, 30 eligible queries were selected and regarded as the initial task queries in our study.

Table 1 shows a set of example search tasks. For each task, the presented verticals on the initial SERP were either relevant or irrelevant. In order to generate the irrelevant vertical results, we use reformulated queries which are modified from the initial task queries collected from practical search logs [1, 8, 28] to retrieve the vertical. Most of the reformulated queries are generalization or specification as shown in Table 1 while the selection criteria is similar to [1] as to reflect the different facets of the query. Since the irrelevant verticals also contained terms from the original queries, the appearance of the irrelevant vertical hyperlinks was similar to the relevant ones. The aim of this is to make the occurrence of irrelevant verticals more natural and reasonable. Base on the interview of the participants, none of them noticed the SERPs with irrelevant verticals had been modified.

In our experimentation, the participants experienced one of the thirty-one experimental conditions (5 vertical types × 3 position × 2 quality conditions of vertical + 1 organic-only result page) for each task. To make sure that all tasks would be completed with equal opportunities in each condition, we used a Graeco-Latin square design [8, 16] of tasks and conditions. Based on this design, we were able to strictly control the effects of tasks and conditions. The task conditions were divided into six groups. Each group consisted of five vertical types, including textual, encyclopedia, image-only, application-download and news. For each specific type, there were six initial pages, which were split by three positions and two levels of quality, corresponding to the tasks. Thus, each participant who performed in one of these six groups would fairly complete all conditions in the

**Table 1: Search Tasks and Manipulated Off-target Queries to Retrieve Verticals**

| Original Query | Off-target Query | Vertical |
|---|---|---|
| Ancient Greek Architectural style | Ancient Greek | Textual |
| Poems on spring rains | Poems on rains | |
| The $9^{th}$ zone (movie) | The $9^{th}$ zone (novel) | Encyclopedia |
| Covering the Sky (novel) | Covering the Sky (game) | |
| Nike basketball shoes | Nike football shoes | Image-only |
| How to cook spiced egg | How to cook omelet | |
| iTunes download | iTools download | Application |
| Renren desktop app download | Weibo desktop app download | |
| Ebola virus mutation | Ebola virus | News |
| Shanghai $3^{rd}$ airport | Shanghai flight | |

experimentation. Furthermore, to eliminate the effects of task order or other possible learning biases [18, 19, 20], the tasks in each group were presented in a random order. After asking the participants whether they had noticed any variations in search performance among those tasks, none of them detected any differences among these conditions.

## 4. EXAMINATION BEHAVIOR ANALYSIS

In the following section, we describe the main findings of the users' examination patterns on the SERP with the presence of the vertical results, compared to the ones with only organic results. Especially, we want to know how users' attention are distributed across different page elements (embedded vertical, surrounding organic results). Specifically, we aim to answer the following research questions (RQs):

- (**RQ1**) Are there any attraction biases in the presence of vertical results? How do the vertical type, embedded position and relevance affect these biases?
- (**RQ2**) Are there any effects in users' attention to organic results around verticals? How do the vertical type, embedded position and relevance affect these biases?
- (**RQ3**) Given different verticals and their presentation styles, what are the examination probabilities of different items within the internal vertical block? What are the common browsing patterns among these items?

Since the SERPs in our experiment are organized in a blended fashion [1, 28], i.e. the results (either organic or vertical) are formed into blocks and ranked from top to bottom, we can simply track users' interactions on different ranked positions. In this work, it is assumed that the SERP generally consists of nine organic results and one vertical result block that is embedded within the organic results. In the following sections, we first investigate user attention in terms of gaze duration on different results (ranked positions) in Section 4.1, followed by a detailed analysis of the gaze duration and movement patterns within vertical blocks in Section 4.2.

## 4.1 Examination Behavior on SERPs with Verticals

First, we present in Table 2 the overall results of how users allocate their attentions on the SERPs based on the different independent variables that we manipulated. The user attention is measured by the fixation distribution we obtained from the eye-tracker. The top part of Table 2 shows the attention distribution of the pure organic result lists and is treated as our baseline for comparison. Note that due to the small fraction of users' attention on the results ranked below position 5 (e.g 6.81% for organic only SERPs), we present the results of position 6-10 as

**Table 2: The user attention (eye fixation distribution) of each ranked position on the SERP when different types and quality of verticals are embedded at various positions on the SERPs. Two-tailed t-test is performed to detect any significant changes against the user attention on organic results only SERPs. Significant results are bold while * and ** represent p < 0.05 and 0.01, respectively. The block color represents the user attention's change of direction compared with organic only SERPs (red denotes increment and blue represents decrement) while the brightness of the color indicates the (normalized) change magnitude.**

| Position | 1 | 2 | 3 | 4 | 5 | 6~10 | Position | 1 | 2 | 3 | 4 | 5 | 6~10 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| No Vertical | | | | | | | | | | | | | |
| | 46.03% | 27.92% | 9.80% | 5.43% | 4.01% | 6.81% | | 46.03% | 27.92% | 9.80% | 5.43% | 4.01% | 6.81% |
| Relevant | | | | | | | Irrelevant | | | | | | |
| Textual | | | | | | | Textual | | | | | | |
| 1 | **61.93%*** | 17.79% | 9.36% | 5.53% | 3.68% | **1.71%*** | 1 | 35.32% | 30.07% | 9.44% | **11.22%*** | 6.55% | 7.39% |
| 3 | 43.39% | 24.90% | **16.13%*** | 3.05% | 8.30% | 4.24% | 3 | 41.24% | 27.35% | 9.99% | 6.73% | 3.48% | 11.20% |
| 5 | 41.15% | 20.96% | 10.25% | **14.93%*** | 4.42% | 8.29% | 5 | 43.26% | 22.07% | 8.77% | **11.61%*** | 7.79% | 6.49% |
| Encyclopedia | | | | | | | Encyclopedia | | | | | | |
| 1 | **61.6%*** | 24.90% | 6.63% | **1.25%*** | 2.21% | 3.41% | 1 | 33.09% | 31.59% | 14.56% | 6.66% | 4.75% | 9.35% |
| 3 | 45.15% | 20.88% | **24.97%*** | 2.49% | 4.00% | 2.51% | 3 | 43.06% | 18.20% | 11.68% | **10.46%*** | 6.15% | 10.46% |
| 5 | 51.61% | **17.1%*** | 9.45% | 7.49% | **10.89%*** | 3.47% | 5 | 54.15% | **17.27%*** | 7.03% | 6.16% | 7.04% | 8.35% |
| Image-only | | | | | | | Image-only | | | | | | |
| 1 | **78.37%*** | **10.48%*** | 8.12% | **1.22%*** | 1.15% | **0.66%*** | 1 | 42.36% | 32.85% | 10.31% | 6.17% | 4.09% | 4.22% |
| 3 | 38.36% | **15.88%*** | **40.5%*** | 3.70% | 0.91% | **0.65%*** | 3 | 37.13% | 18.22% | **21.89%*** | 8.14% | 5.34% | 9.28% |
| 5 | **20.47%*** | **16.41%*** | 12.88% | **10.91%*** | **34.4%*** | 4.92% | 5 | **27.67%*** | **12.66%*** | 15.34% | **11.43%*** | **23.42%*** | 9.48% |
| Application-download | | | | | | | Application-download | | | | | | |
| 1 | **81.36%*** | **5.7%*** | **2.57%*** | 2.49% | 5.52% | 2.36% | 1 | 41.99% | 28.48% | 12.98% | 5.25% | 3.85% | 7.45% |
| 3 | **21.62%*** | 19.49% | **50.12%*** | 2.10% | 4.18% | 2.48% | 3 | 35.86% | **17.2%*** | **35.28%*** | 6.59% | 2.85% | 2.22% |
| 5 | 43.62% | 17.83% | 4.64% | 7.34% | **25.06%*** | **1.5%*** | 5 | 46.48% | 20.91% | 13.27% | 4.37% | **11.03%*** | 3.95% |
| News | | | | | | | News | | | | | | |
| 1 | 50.55% | **10.99%*** | 10.48% | **10.28%*** | 3.73% | **13.97%*** | 1 | 42.85% | 24.88% | 14.15% | 3.88% | 3.66% | 10.59% |
| 3 | 33.19% | 18.94% | **29.37%*** | 5.40% | 6.80% | 6.29% | 3 | 34.28% | 33.01% | **17.35%*** | 6.20% | 2.75% | 6.40% |
| 5 | 44.17% | 19.09% | 10.22% | 9.47% | 7.97% | 9.08% | 5 | 34.91% | 30.12% | 14.61% | 5.78% | 4.49% | 10.08% |

an aggregation. Not surprisingly, aligned with previous eye-tracking results [9, 28], the users' attention on the organic-only results decreases as the ranked position increased on the SERP.

We then present the attention distributions of various verticals that are embedded in different positions (1, 3, 5). In Table 2, the left and right side of the table respectively present the results of when the vertical is relevant and irrelevant. In Table 2, we can observe that when there is no vertical (organic-only SERP), 46.03% and 27.92% of users' fixations are focused on position 1 and 2, respectively. If there is a relevant textual vertical embedded at position 1 on the SERP, the fixation on position 1 (vertical) increases to 61.93%, meanwhile the one of the position 2 (web) decreases to 17.79%. This demonstrates that the vertical embedded in the first result of the SERP significantly (*p-value*=0.05) attracts more user attention compared to a presented organic web result. Similarly, by examining and comparing the attention distribution in Table 2, we can observe a number of interesting findings, as discussed below.

### 4.1.1 Attraction Effect

Firstly, we focus on users' attention allocation on the vertical results themselves (to answer RQ1), i.e. tracking solely the user attention on the corresponding vertical embedded position: 1, 3, 5. As shown in Table 2, we find a strong attraction bias towards vertical results once it is presented (compared to the corresponding one on the organic-only SERP). This attraction effect is strong for almost all different types of verticals at all embedded positions (1, 3, 5) when the vertical results are relevant (left side of Table 2). We can observe that there are significant increases over the organic-only SERP except for positions 1 and 5 of news vertical. This attraction bias is highly significant (*p-value*<0.01) for image-only and application-download verticals for all positions, position 3 and 5 for encyclopedia and position 3 for news verticals. From the perspective of vertical embedding positions, the attraction bias is strongest in the 3rd ranking position where all verticals obtain statistically significant increases over the organic-only SERP in terms of attention.

When the vertical is irrelevant (right side of Table 2), we can observe that there is still a strong bias towards verticals on positions 3 and 5, especially for image-only and application-download verticals. Although as expected, all the attraction bias is diminished compared to the cases where verticals are relevant. The results above are interesting but not surprising, as users can easily distinguish a vertical result block from organic results (e.g. image-only, application-download). When the vertical is relevant (potentially the vertical results are what the users are looking for), the users tend to pay more attention to the vertical results. Note that although this attraction bias effect has also been noticed in previous eye-tracking based studies [28], we investigate more verticals with various presentation styles. We also study how the vertical relevance (or namely orientation in [27, 29]) affects
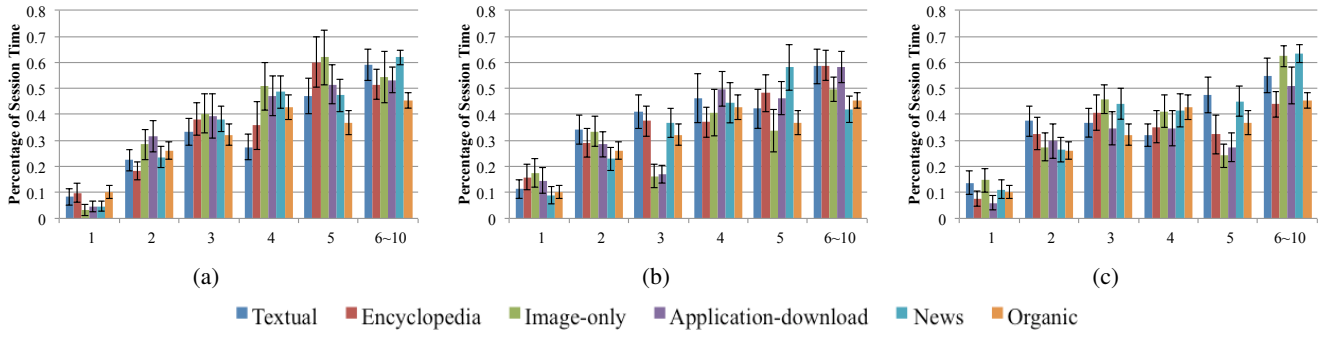
Figure 2: Mean time of arrival at each position when different verticals are placed at position 1, 3, 5 respectively.

the user's interaction and reaffirms previous results through eye-tracking (rather than the click-through analysis).

To further analyze this effect in terms of whether users are attracted by the vertical results and tend to examine them first (with respect to the temporal sequence), we analyze the arrival time distribution of different ranked positions when the various vertical results are embedded in different positions. The results are shown in Figure 2. The x-axis is the ranked positions while the y-axis is the first arrival of eye fixation for a given position (normalized by the whole session length due to user variabilities [27]). We can observe from Figure 2(a) that when there is a vertical shown at position 1 with a strong visual bias in the snippets (e.g. image-only, application and news verticals), users' arrival time to the vertical hugely decreases, i.e. the users recognize the vertical immediately. Due to this visual bias, users also tend to delay their reading of the organic results presented at the other positions compared to organic-only SERP. From Figure 3(b), we can see that users are still attracted to image-only and application-download vertical presented at position 3, but this is not the case for news, encyclopedia and textual verticals. Interestingly, for the news vertical, the users' first arrival time at position 3 is quite different from when shown in the first position (weaker attraction bias). This can be explained by the fact that information in the news vertical may be already covered in organic results with similar content. Therefore, there is no need for users to pay so much attention to news verticals. Figure 3(c) shows similar results while image-only and application-download verticals attract much attention on position 5.

### 4.1.2 Cut-off and Spill-over Effect

We are also interested in users' attention changes in the organic results after users' examination of the vertical results (to answer **RQ2**).

To study this, firstly, we plot in Figure 3 the mean percentages of eye fixation durations on organic results after users have examined the verticals. We present both cases when the vertical presented is relevant or irrelevant. From this figure, we can see that after users examine the vertical results, the users tend to pay less attention to organic results if the vertical is relevant and allocate more attention to organic ones if the vertical is irrelevant. It is interesting that the user behavior has changed when the relevance of the vertical varies.

Now we look into these two different cases in more details. When the vertical is relevant, we can observe in Table 2 that compared to the organic only results, the user attention on the results ranked below the vertical in general decreased. This is especially true for the image-only and application-download
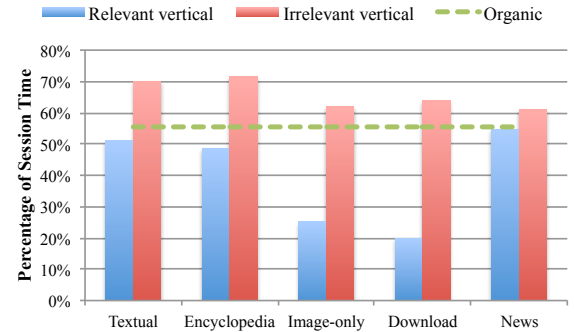


Figure 3: Mean percentage of fixation durations on other surrounding organic results after the users have examined vertical results with different types and qualities (relevance or not). For the organic case (baseline), it is after the user have examined the organic results shown at the corresponding position as same as the vertical position (at Position 1).

verticals, although not all of the engagement differences are significant.

To study how users pay attention to the organic results ranked below the vertical results after they have examined the relevant vertical results, we present the percentage of fixation duration on these lowly ranked organic results and their normalized differences with the organic-only SERPs in Table 3. We can see that after users have viewed the vertical results, they tend to decrease their visual attention on the organic results which are below these verticals, especially when the image-only and application-download verticals was embedded at position 3. Overall, we refer to this effect of paying less attention to the organic results ranked below the relevant vertical as the *cut-off effect*. From Table 3, we can observe that the *cut-off effect* not only depends on the types of verticals, but is also influenced by the verticals' positions. For example, when the verticals are at position 5, the decrease of attention on organic results ranked below the verticals is relatively limited, even for image-only and application-download verticals.

However, when the vertical is irrelevant, we can see in Table 2 that compared to organic only SERP, the users consistently pay more attention to the results after the vertical's presence but less to the results prior to the vertical embedded position. This is *spill-over effect* and it has also been observed in previous work [1] although our results indicate that the spill-over effect happens more often to those documents posterior to the irrelevant verticals.

198

**Table 3: The percentage of fixation duration on organic results ranked below the vertical (the same position for the organic case). This is after the users have examined the relevant vertical results that vary in different embedded position (3 or 5). The difference (Diff) is calculated as (Vertical-Organic)/Organic and two-tailed t-test is performed for significance.**

| Relevant Vertical | Textual | Encyclo-pedia | Image-only | Application-download | News |
|---|---|---|---|---|---|
| | Position = 3 | | | | |
| Organic | 34.61% | | | | |
| Vertical | 30.13% | 16.70% | 8.44% | 13.04% | 22.61% |
| Diff | -12.95% | **-51.74%*** | **-75.62%**** | **-62.32%**** | -34.68% |
| | Position = 5 | | | | |
| Organic | 25.27% | | | | |
| Vertical | 26.30% | 19.27% | 10.33% | 6.21% | 38.69% |
| Diff | 4.09% | -23.76% | **-59.10%*** | **-75.44%*** | 53.09% |

**Table 4: The percentage of fixation duration on organic results surrounding the vertical. This is after the users have examined the vertical results that vary in different embedded position (3 or 5) and vertical relevance (relevant or irrelevant). The difference (Diff) is calculated as (Irrelevant-Relevant)/Relevant and two-tailed t-test is performed for significance.**

| | Textual | Encyclo-pedia | Image-only | Application-download | News |
|---|---|---|---|---|---|
| | Position = 3 | | | | |
| Relevant | 61.49% | 59.90% | 45.14% | 52.37% | 54.86% |
| Irrelevant | 86.72% | 87.74% | 76.81% | 63.14% | 64.07% |
| Diff | **41.04%*** | **46.47%**** | **70.18%**** | 20.57% | 16.80% |
| | Position = 5 | | | | |
| Relevant | 80.97% | 59.56% | 46.58% | 49.64% | 69.81% |
| Irrelevant | 79.92% | 91.57% | 68.93% | 80.01% | 85.93% |
| Diff | -1.30% | **53.75%**** | 47.99% | **61.16%**** | **23.09%*** |

To further examine this spill-over effect, we also present the user attention of organic results with the relevant and irrelevant vertical results after the users have examined the verticals in Table 4. By comparing this organic result examination distribution for both relevant and irrelevant verticals, we found that it is the case that the users pay more attention to the organic results for irrelevant vertical and there are significantly differences found in encyclopedia and image-only vertical (*p-value*<0.01) between the relevant and irrelevant verticals. We also observe significant differences at position 5 when the types of verticals are encyclopedia and application-download (*p-value*<0.01). This demonstrates that this spill-over effect varies according to the relevance of verticals and can cause more visual attention on the surrounding organic results, especially when the vertical is irrelevant.

### 4.1.3 Summary

So far, we have extensively analyzed how users examine the vertical and the surrounding organic results on verticals of different types, embedded positions and relevance. To conclude, the findings (for RQ1 and RQ2) of user attention (measured by the gaze duration) on SERPs with verticals in the previous two subsections can be summarized as below:

- *attraction bias*: there is a strong bias towards more attention on the vertical results when the vertical results are relevant or even irrelevant but visually appealing;
- *cut-off effect*: when the vertical results are relevant, there is a strong bias towards paying less attention to the organic results posterior to the relevant vertical;
- *spill-over effect*: irrelevant verticals would increase the attention for organic results.

## 4.2 Examination Behavior within Vertical Blocks

Although we have gained lots of insights on how users examine results on the SERP with vertical results, however, little is known about how users examine the various items (components) within vertical blocks.

In this section, we aim to understand further the user examination patterns within different vertical blocks (to answer **RQ3**). Specifically we study: (a). the examination probabilities of viewing the various components within the verticals with various presentation styles; (b). the most common examination patterns, especially those reflect how users start and end their vertical examination. We focus on these two research questions for the following reasons:

- A good understanding of the examination probability can help better estimate the relevance of the item within the vertical, which could be useful in vertical-aware click models [9]. So far most of the research only considers using click/skip on the vertical level while utilizing the examination and click of the items within the vertical can be a more fine-grained feedback.
- Understanding better on the browsing patterns under different presentation styles could be useful in studying the decision process of how user percieve the relevance of the vertical. The sequential information could also help in deciding how to best present the ranked items in order to fit the most relevant one to the first user examined zone.

With the assistance of eye-tracking devices, Figure 4 presents the results of the internal examination process of users. Figure 4(a) shows the examination distributions (measured by fixation) of different components (items, vertical title, etc.) within the vertical block, while (b) shows the vertical block presentation style for the five verticals. Figure 4(c) shows the detailed information on how users examine within the vertical block, specifically on the top 5 most common examination patterns and their corresponding distributions of how they start the examination and where they end their examination. For example, in Figure 4(c) (textual vertical), we can observe that the most popular examination pattern is "(1, END)" which means that the users directly examined zone 1 (i.e. item1) and then end their vertical block examination. Similarly, the second popular examination pattern "(0, 1, END)" indicates the users first examined zone 0 (vertical title) and then zone 1 (item1) before ending their examination.

With respect to Figure 4(a), we can observe several interesting trends among different vertical types. Firstly, we find that a substantial amount of attention (around 30%) has been paid to the vertical title (besides news vertical), which suggests that users pay attention to the vertical type when examining/judging the relevance of the vertical block. Secondly, the users may pay more attention to the items that contain substantial information (e.g. detailed textual snippets) while the examination probabilities for different items within the vertical generally follow the trend that the higher ranked items receive higher attention. Additionally, it is interesting to observe that users tend to allocate their visual attention almost equally on these similar components and the position of these vertical components do not affect users' internal
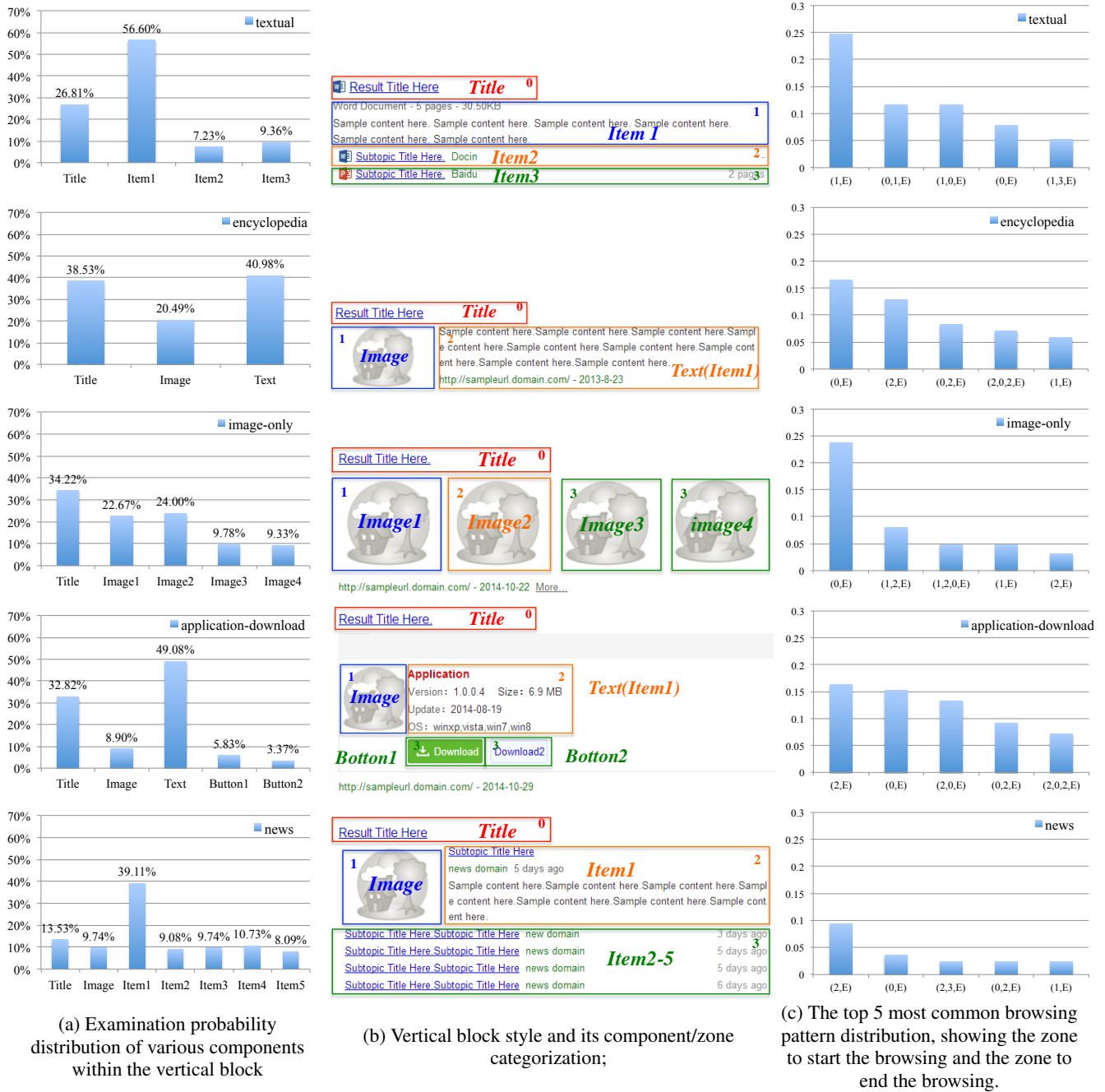
(a) Examination probability distribution of various components within the vertical block

(b) Vertical block style and its component/zone categorization;

(c) The top 5 most common browsing pattern distribution, showing the zone to start the browsing and the zone to end the browsing.

Figure 4: **Examination of various vertical blocks: textual, encyclopedia, image, application download and news. Note that the different border color represents different zones in the browsing pattern distribution (red, blue, orange, green denotes zone 0, 1, 2, 3, respectively) for the ease of presentation. These borders of the zones are not displayed during our experimentation.**

examination. For example, the percentages of fixation durations are almost similar (around 9%) from item 2 to item 5 in the news vertical.

We also find different attraction effects between text and image in the vertical blocks. By comparing the eye fixations on the text and image of encyclopedia, application-download and news verticals, although image sometimes (e.g. 20% for image-only) is examined, interestingly, the bias of this when examining in the vertical block is not as high as we expected. This suggests that although image in the vertical is visually appealing and could trigger strong attraction bias when the users examine the SERP, the impact of this visual attractiveness is not high after the users has decided to further examine the vertical block. Besides this, the influence of component size on the users' internal examination is also shown in Figure 4(a). We can see that, not surprisingly, the larger components might receive more visual attention and users tend to examine the components which contain substantial information (e.g. detailed textual snippets) in the vertical.

Now we discuss the examination pattern for each vertical. For the textual vertical, we found that most of the attention is paid to the first item while generally the users start or end their vertical examination using the vertical title or the first item (item1). For the encyclopedia vertical, we found that the users, similarly, tend to start and end their sessions by looking at the vertical title or the textual document while seldom using the images as the starting/ending point to examine these vertical blocks. For the image-only vertical, it is interesting to find that the second image item receives more attention (see Figure 4(a)) and is examined more often at the end of the session than the first image item. That suggests that users tend to end their internal examination within image-only vertical after they viewed the second image items. For the application and news vertical, we found similar results to the encyclopedia vertical that users mainly examines item1 and vertical title to decide vertical relevance and the image visual bias is also not high in the vertical block examination stage.

To summarize, in this section, we study the examination patterns within the vertical block and found the following findings:

- The users tend to examine the vertical type (title) and the detailed snippets of the first item to judge the relevance of the vertical and tend to skip the further results if they expect it to be irrelevant.
- After the users decide to examine the vertical block, the image shown within the vertical has relatively weak effect towards users' examination.
- The position of homogeneous components in verticals might weakly affect users' visual attention. Users tend to treat the similar components equally when they examine these items within the vertical blocks.
- The component size has a strong effect on users' internal examination. When the components occupy more area of vertical block and contain detailed information, users are more likely to examine this components and allocate more visual attention on it.

## 4.3 Summary

In this section, we summarize the main findings and limitations of our study.

### 4.3.1 Behavior on SERPs

Our first research question (RQ1) investigates whether attraction biases exist in the presence of vertical result. Focusing on the fixation distribution (see Section 4.1), we observe that not all verticals maintain strong attraction biases when the SERPs contain verticals. Essentially, the attraction effect is influenced by the types of verticals, while the vertical quality (relevant or not) does not have a huge impact. Table 5 presents the overview of different effects according to vertical type and quality. From Table 5, it is clear that users are more likely to be attracted by the image-only and application-download verticals than the news ones. Interestingly, news verticals, whose user performance differs much from the others, appear not to capture more visual attention than the organic results (even news contain image in the vertical block). Based on our eye-tracking analysis, we find that users may ignore the image within the news vertical block and treat the news items equally with the other organic results.

We are also interested in whether verticals affect users' attention on organic results (RQ2). Based on the analysis of user interactions after examining the vertical results, we find users may significantly change their behavior on organic results when they notice the presence of vertical. It is interesting to show that users tend to reduce their attention on the organic results ranked below

**Table 5: The overview of three user behavior effect on the SERPs. "●" represents the significant results with the corresponding effect, and "○" represents the weak effect (exists but is not significant) of the corresponding vertical type. "−" means the corresponding effect did not perform in this vertical.**

|  | Attraction Effect | | Cut-off Effect | | Spill-over Effect |
|---|---|---|---|---|---|
|  | Relevant | Irrelevant | Relevant | Irrelevant |  |
| Textual | ● | - | - | - | ○ |
| Encyclo-pedia | ● | ○ | ○ | - | ● |
| Image-only | ● | ● | ● | - | ● |
| Application-download | ● | ● | ● | - | ● |
| News | ○ | ○ | - | - | ○ |

the relevant vertical, especially for image-only and application-download. The results of our study reveal that this *cut-off effect* widely exists in the federated search behavior. However, as shown in Table 5, we can observe that this is not the case when the irrelevant vertical results are presented. Besides this, we also focus on the spill-over effect, which has been observed in prior work. We further demonstrate that users would pay more attention on organic results when the irrelevant vertical is placed on the SERPs.

### 4.3.2 Internal Examination within Vertical

Our third research question (RQ3) focus on the internal examination process within the vertical blocks. In our experiment, it is interesting to observe that different vertical types would also lead to a wide variety of internal examination behavior. We find that users perfer to examine the vertical title and the detail snippets of the first item to judge the relevance and skip the further results. Another interesting finding is that images do not perform a strong attraction effect when users has decided to examine the vertical blocks. In our experiment, we also find the size of components may also affect users' visual attention.

### 4.3.3 Limitations

There are several potential limitations in our study. Firstly, we mainly focused on informational and transactional search tasks and did not explore tasks with navigational search intent [6, 20]. Because the first result for the navigational search is often the destination site, it is difficult to distinguish the effects of verticals from the task type bias. Secondly, to concentrate on the influence of verticals and create a controlled lab experiment, we did not take multiple verticals into account. In real life, more than one vertical might be presented on SERPs. This might also affect different user behavior.

## 5. CONCLUSIONS AND FUTURE WORK

As more vertical results are adopted in the Web search, it is necessary to better understand the influence of vertical results on user decision strategies. In this paper, we present an in-depth study using eye-tracking techniques to explore user examination behavior on federated search. Based on the data collected from a controlled experimentation, we systematically analyze user attention both at page level and result level. Our results suggest three interesting effect in federated search. Firstly, we find attraction effect of verticals, which has been observed in existing studies [9, 28], is affected by vertical types and users' search intention. While image-only and application-download verticals

has the strongest attraction effect, it is interesting to note that news verticals seem not to attract more users attention than organic results. Secondly, we demonstrate, for the first time, that user may drastically reduce their visual attention on organic results, after examining the relevant verticals. In other words, users' examination on organic may be "cut off" by the placement of vertical results. In our experimentation, the SERPs with application-download verticals perform the most significant *cut-off effect*, followed by the image-only. Then, we confirm the *spill-over effect* also has an influence on users' examination behavior. To take it a step further, we find that the spill-over effect has some special "orientations" in federated search, which means that users prefer to pay more attention on the organic results, when the verticals are irrelevant or off topic.

The internal examination of verticals is another key concern in our study. As most of prior work focused on the vertical effect at page level, this is, to our knowledge, the first quantitative investigation on the internal examination of verticals. Interestingly, our findings show there is a shrap difference in users' attention for the different types of vertical. Interesting directions for future work involve extending this work and developing the popular adopted evaluation metrics based our findings so that they can better correlate with the user preferences on the SERPs with verticals. Moreover, predicting and modeling the users behavior with these three effects on federated search is another interesting challenge in our future work.

## 6. ACKNOWLEDGMENTS

## 7. REFERENCES

[1] J. Arguello and R. Capra. The effects of vertical rank and border on aggregated search coherence and search behavior. In *Proceedings of the 23rd ACM International Conference on Conference on Information and Knowledge Management*, pages 539–548. ACM, 2014.

[2] J. Arguello, F. Diaz, and J. Callan. Learning to aggregate vertical results into web search results. In *Proceedings of the 20th ACM international conference on Information and knowledge management*, pages 201–210. ACM, 2011.

[3] J. Arguello, F. Diaz, J. Callan, and J.-F. Crespo. Sources of evidence for vertical selection. In *Proceedings of the 32nd international ACM SIGIR conference on Research and development in information retrieval*, pages 315–322. ACM, 2009.

[4] A. Aula, P. Majaranta, and K.-J. Räihä. Eye-tracking reveals the personal styles for search result evaluation. In *Human-Computer Interaction-INTERACT 2005*, pages 1058–1061. Springer, 2005.

[5] J. Bar-Ilan, K. Keenoy, M. Levene, and E. Yaari. Presentation bias is significant in determining user preference for search results - a user study. *Journal of the American Society for Information Science and Technology*, 60(1):135–149, 2009.

[6] A. Broder. A taxonomy of web search. In *ACM Sigir forum*, volume 36, pages 3–10. ACM, 2002.

[7] G. Buscher, A. Dengel, and L. van Elst. Eye movements as implicit relevance feedback. In *CHI'08 extended abstracts on Human factors in computing systems*, pages 2991–2996. ACM, 2008.

[8] G. Buscher, S. T. Dumais, and E. Cutrell. The good, the bad, and the random: an eye-tracking study of ad quality in web search. In *Proceedings of the 33rd international ACM SIGIR conference on Research and development in information retrieval*, pages 42–49. ACM, 2010.

[9] D. Chen, W. Chen, H. Wang, Z. Chen, and Q. Yang. Beyond ten blue links: enabling user click modeling in federated web search. In *Proceedings of the fifth ACM international conference on Web search and data mining*, pages 463–472. ACM, 2012.

[10] E. Cutrell and Z. 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*, pages 407–416. ACM, 2007.

[11] F. Diaz, R. White, G. Buscher, and D. Liebling. Robust models of mouse movement on dynamic web search results pages. In *Proceedings of the 22nd ACM international conference on Conference on information & knowledge management*, pages 1451–1460. ACM, 2013.

[12] S. T. Dumais, G. Buscher, and E. Cutrell. Individual differences in gaze patterns for web search. In *Proceedings of the third symposium on Information interaction in context*, pages 185–194. ACM, 2010.

[13] J. Fernquist and E. H. Chi. Perception and understanding of social annotations in web search. In *Proceedings of the 22nd international conference on World Wide Web*, pages 403–412. International World Wide Web Conferences Steering Committee, 2013.

[14] L. A. Granka, T. Joachims, and G. Gay. Eye-tracking analysis of user behavior in www search. In *Proceedings of the 27th annual international ACM SIGIR conference on Research and development in information retrieval*, pages 478–479. ACM, 2004.

[15] Q. Guo and E. Agichtein. Towards predicting web searcher gaze position from mouse movements. In *CHI'10 Extended Abstracts on Human Factors in Computing Systems*, pages 3601–3606. ACM, 2010.

[16] K. Hofmann, B. Mitra, F. Radlinski, and M. Shokouhi. An eye-tracking study of user interactions with query auto completion. In *Proceedings of the 23rd ACM International Conference on Conference on Information and Knowledge Management*, pages 549–558. ACM, 2014.

[17] J. Huang, R. White, and G. Buscher. User see, user point: gaze and cursor alignment in web search. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, pages 1341–1350. ACM, 2012.

[18] T. Joachims, L. Granka, B. Pan, H. Hembrooke, and G. Gay. Accurately interpreting clickthrough data as implicit feedback. In *Proceedings of the 28th annual international ACM SIGIR conference on Research and development in information retrieval*, pages 154–161. ACM, 2005.

[19] T. Joachims, L. Granka, B. Pan, H. Hembrooke, F. Radlinski, and G. Gay. Evaluating the accuracy of implicit feedback from clicks and query reformulations in web search. *ACM Transactions on Information Systems (TOIS)*, 25(2):7, 2007.

[20] D. Lagun, C.-H. Hsieh, D. Webster, and V. Navalpakkam. Towards better measurement of attention and satisfaction in mobile search. In *Proceedings of the 37th international ACM SIGIR conference on Research & development in information retrieval*, pages 113–122. ACM, 2014.

[21] X. Li, Y.-Y. Wang, and A. Acero. Learning query intent from regularized click graphs. In *Proceedings of the 31st annual international ACM SIGIR conference on Research and development in information retrieval*, pages 339–346. ACM, 2008.

[22] Y. Liu, C. Wang, K. Zhou, J. Nie, M. Zhang, and S. Ma. From skimming to reading: A two-stage examination model for web search. In *Proceedings of the 23rd ACM International Conference on Conference on Information and Knowledge Management*, pages 849–858. ACM, 2014.

[23] V. Navalpakkam, L. Jentzsch, R. Sayres, S. Ravi, A. Ahmed, and A. Smola. Measurement and modeling of eye-mouse behavior in the presence of nonlinear page layouts. In *Proceedings of the 22nd international conference on World Wide Web*, pages 953–964. International World Wide Web Conferences Steering Committee, 2013.

[24] M. Richardson, E. Dominowska, and R. Ragno. Predicting clicks: estimating the click-through rate for new ads. In *Proceedings of the 16th international conference on World Wide Web*, pages 521–530. ACM, 2007.

[25] K. Rodden and X. Fu. Exploring how mouse movements relate to eye movements on web search results pages. *Web Information Seeking and Interaction*, pages 29–32, 2007.

[26] K. Rodden, X. Fu, A. Aula, and I. Spiro. Eye-mouse coordination patterns on web search results pages. In *CHI'08 Extended Abstracts on Human Factors in Computing Systems*, pages 2997–3002. ACM, 2008.

[27] S. Sushmita, H. Joho, M. Lalmas, and R. Villa. Factors affecting click-through behavior in aggregated search interfaces. In *Proceedings of the 19th ACM international conference on Information and knowledge management*, pages 519–528. ACM, 2010.

[28] C. Wang, Y. Liu, M. Zhang, S. Ma, M. Zheng, J. Qian, and K. Zhang. Incorporating vertical results into search click models. In *Proceedings of the 36th international ACM SIGIR conference on Research and development in information retrieval*, pages 503–512. ACM, 2013.

[29] K. Zhou, R. Cummins, M. Lalmas, and J. M. Jose. Evaluating aggregated search pages. In *Proceedings of the 35th international ACM SIGIR conference on Research and development in information retrieval*, pages 115–124. ACM, 2012.

[30] K. Zhou, R. Cummins, M. Lalmas, and J. M. Jose. Which vertical search engines are relevant? In *Proceedings of the 22nd international conference on World Wide Web*, pages 1557–1568. International World Wide Web Conferences Steering Committee, 2013.