

The Infinite Index: Information Retrieval on Generative Text-To-Image Models

NIKLAS DECKERS, Leipzig University and ScaDS.AI, Germany

MAIK FRÖBE, Friedrich-Schiller-Universität Jena, Germany

JOHANNES KIESEL, Bauhaus-Universität Weimar, Germany

GIANLUCA PANDOLFO, Bauhaus-Universität Weimar, Germany

CHRISTOPHER SCHRÖDER, Leipzig University, Germany

BENNO STEIN, Bauhaus-Universität Weimar, Germany

MARTIN POTTHAST, Leipzig University and ScaDS.AI, Germany

The text-to-image model Stable Diffusion has recently become very popular. Only weeks after its open source release, millions are experimenting with image generation. This is due to its ease of use, since all it takes is a brief description of the desired image to “prompt” the generative model. Rarely do the images generated for a new prompt immediately meet the user’s expectations. Usually, an iterative refinement of the prompt (“prompt engineering”) is necessary for satisfying images. As a new perspective, we recast image prompt engineering as interactive image retrieval—on an “infinite index”. Thereby, a prompt corresponds to a query and prompt engineering to query refinement. Selected image–prompt pairs allow direct relevance feedback, as the model can modify an image for the refined prompt. This is a form of one-sided interactive retrieval, where the initiative is on the user side, whereas the server side remains stateless. In light of an extensive literature review, we develop these parallels in detail and apply the findings to a case study of a creative search task on such a model. We note that the uncertainty in searching an infinite index is virtually never-ending. We also discuss future research opportunities related to retrieval models specialized for generative models and interactive generative image retrieval. The application of IR technology, such as query reformulation and relevance feedback, will contribute to improved workflows when using generative models, while the notion of an infinite index raises new challenges in IR research.

ACM Reference Format:

Niklas Deckers, Maik Fröbe, Johannes Kiesel, Gianluca Pandolfo, Christopher Schröder, Benno Stein, and Martin Potthast. 2022. The Infinite Index: Information Retrieval on Generative Text-To-Image Models. 1, 1 (December 2022), 19 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1 INTRODUCTION

Conditional generative models allow the generation of a desired output based on a user-specified condition. For generative text-to-image models such as DALL-E [57] or Stable Diffusion [59], this means that the model generates images that match a given text description, the “prompt.” For a user, the prompt is the primary means of controlling

Authors’ addresses: Niklas Deckers, niklas.deckers@uni-leipzig.de, Leipzig University and ScaDS.AI, Center for Scalable Data Analytics and Artificial Intelligence, Dresden/Leipzig, Germany; Maik Fröbe, maik.froebe@uni-jena.de, Friedrich-Schiller-Universität Jena, Jena, Germany; Johannes Kiesel, johannes.kiesel@uni-weimar.de, Bauhaus-Universität Weimar, Weimar, Germany; Gianluca Pandolfo, gianluca.pandolfo@uni-weimar.de, Bauhaus-Universität Weimar, Weimar, Germany; Christopher Schröder, christopher.schroeder@uni-leipzig.de, Leipzig University, Leipzig, Germany; Benno Stein, benno.stein@uni-weimar.de, Bauhaus-Universität Weimar, Weimar, Germany; Martin Potthast, martin.pothast@uni-leipzig.de, Leipzig University and ScaDS.AI, Center for Scalable Data Analytics and Artificial Intelligence, Dresden/Leipzig, Germany.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2022 Association for Computing Machinery.

Manuscript submitted to ACM

Table 1. Comparison of the most relevant text-to-image models. The * refers to replicated resources.

Text-to-image model		Training data		Resource publication			
Name	Parameters	Size	Source	Code	Data	Model	Month / Year
DALL·E [57]	12 B	N/A	Custom web crawl ¹	✓*	✗	✓*	01 / 2021
DALL·E 2 [56]	3.5 B	N/A	Custom web crawl, licensed sources ²	✗	✗	✗	07 / 2022
Imagen [61]	4.6 B	400 M [66]	Common Crawl	✓*	✓	✓*	05 / 2022
Midjourney [64]	N/A	N/A	N/A	✗	✗	✗	07 / 2022
Stable Diffusion [59]	890 M	400 M [66]	Common Crawl	✓	✓	✓	08 / 2022

¹ https://github.com/openai/DALL-E/blob/master/model_card.md#training-data² <https://github.com/openai/dalle-2-preview/blob/main/system-card.md#model-training-data>

the generated image. If an ad hoc prompt does not produce a satisfactory result, the user usually interacts with the model by changing the prompt until they get a satisfactory result (unless they give up after a few tries). Since such trial and error is often necessary to achieve satisfactory results, users have begun to exchange ideas and techniques commonly referred to as “prompt engineering.” But aside from referring to and adapting examples from others, it’s often not obvious how to change the prompt to steer image generation in a particular direction. As a new perspective on the use of conditional generative models, we conceptualize them as a novel type of index where a prompt takes the role of a query representing a user’s information need as in text-based image retrieval. Prompt engineering can then be considered a form of interactive text-based image retrieval, in which a user interacts with the model by modifying their prompt as if refining their query in an image search engine to “find” an image that satisfies their information needs.

Apart from the outlined parallels, retrieval on the basis of generative models poses a number of new challenges: At present, the initiative lies solely with the user, without support from the model as a “retrieval system”. In this regard, the “system side” lacks a retrieval model as an intermediary between the user and the model as an index, tailored to help users obtain satisfactory images fast(er), if not ad hoc. For example, the user’s manual refinement of prompts as queries is not supported by system-side log analysis and query expansion, and there is also no operationalization of the concept of relevance, especially when many (different) images are generated. Perhaps the most striking difference from traditional image retrieval is that when generative text-image models are used as an index, new images are generated rather than existing images retrieved. Instead of retrieving relevant images for a query from a finite set of indexed images, a result is returned for every conceivable prompt as a query. This includes prompts for which a traditional image search engine would return no results. Moreover, the number of different images that can be generated per prompt as a query is not conceptually limited, but only by the available computational capacity for model inference. A generative model is therefore effectively an “infinite index.”

Our contribution is to explore this perspective on generative models as indexes in four ways, focusing on text-image generation: (1) We conduct a literature survey on image generation, text-based image retrieval, retrieval for creative tasks, and interactive retrieval (Section 2). (2) We conceptualize generative text-to-image models as an index integrated into a retrieval system: from the user perspective, the query language and interaction methods are presented, and from the system perspective, retrieval technologies capable of supporting retrieval are examined. Requirements for the evaluation of retrieval systems based on generative models are also presented (Section 3). (3) Based on these findings, we conduct a case study on the generation of images for creative tasks in game design that highlights several issues related to the currently available technology (Section 4). (4) Finally, based on the insights gained, an active learning approach to interactive retrieval to guide image generation using generative models is discussed (Section 5).

2 BACKGROUND AND SURVEY

We present related work on generative models and IR to put current developments in context with established concepts.

2.1 Image Generation

In the computer vision field of image synthesis, promising results have been achieved with GANs [6, 23], which allow to generate images from the distribution of given training images. Autoregressive transformer models [57, 58] have been successful in high-resolution image synthesis. Recent developments in diffusion models [68] have shown their ability to outperform the traditional models like GANs in image synthesis [16]. Another major development step was the possibility of conditioning the generated images on text [59]. This forms the foundation of text-to-image models, which are often trained on datasets of text-image pairs [66]. Table 1 gives an overview of the most important recent models. The most recent breakthrough was achieved with the public availability of the model weights of Stable Diffusion [59], which forms a catalyst for many new applications.

2.2 Image Retrieval

While text-to-image generation models are rather new, image retrieval has a long history in research. The literature distinguishes two cases: in content-based image retrieval the user provides an image as query, whereas in text-based image retrieval the user provides a textual query. Content-based image retrieval systems aim to bridge the semantic gap between the high-level meaning of images and their low-level visual features through sophisticated image representations [44]. Having represented and indexed images, the representation of the query image is then used for similarity-based searching and ranking. Text-based image retrieval in the past often focused on retrieval based on image metadata and tags, which is why it is sometimes also referred to as annotation-based, concept-based, or keyword-based image retrieval. However, some approaches also generate textual representations for un-annotated images, for example using character recognition [70], by mapping between text and images using image clusters with known textual representations [46], or by using image captioning methods [30].

Some studies investigated search interactions of users with a text-based image retrieval system. Choi [11] analyzed search logs collected from 29 college students, finding that the participants more often modified their textual queries than that they started a completely new one. Hollink et al. [27] studied the image search behavior of news professionals, showing that these often modified their query by following semantic relationships of query terms, like first looking for images on a person and then for images on their spouse.

Cho et al. [10] took a closer look at why people search for images. In their survey of 69 papers they identified seven information need categories: (1) searching for entertainment, (2) searching for an illustration (explanation or clarification of details, e.g., creating presentation slides or preparing study material), (3) searching for images for aesthetic appreciation (e.g., for desktop backgrounds), (4) searching for knowledge construction (four sub-categories: information processing, information dissemination, learning, and ideation), (5) searching for something to engage others (e.g., to grab audiences' attention), (6) searching for images to inspire themselves, and (7) searching for images for social interactions (e.g., images to trigger emotions). Moreover, they found seven categories of obstacles that could interfere with a user's ability to find the images they seek: (1) semantic problems, i.e., related to employed terminology, (2) content-based issues, i.e., related to describing content of images, (3) technical limitations of systems, (4) issues of aboutness or lacking relevance of retrieved images, (5) problems of inclusivity with regard to cultural or linguistic

aspects of the user, (6) problems of lacking skills in handling search technology, and (7) problems of cognitive overload. As we discuss in Section 3, most of these needs and problems are also relevant for text-to-image generation models.

2.3 Retrieval for Creative Tasks

The use of text-to-image generation models has sparked discussions especially with regards to its artistic or creative uses, thus raising the question whether there are parallels between such use and the literature on search for creative tasks. Interestingly, text-to-image generation models have been able to quickly build communities around their use, not only with regard to tools, but also around sharing prompts for generating images via community feeds and channels.¹ This is in line with artists forming creative communities also in other types of art [24]. On the other hand, such a strong community-building is somewhat surprising as arts and crafts hobbyists are reported to be usually less reliant on human sources [41].

Several studies already analyzed user behavior and goals in creative tasks specifically. Chavula et al. [9] investigated the information behavior for 15 graduate students on creative web search tasks using questionnaires and the think-aloud method. They identified four creative thinking processes between which participants switched back-and-forth: planning for creative search tasks (i.e., deciding on a vague idea), searching for new ideas, synthesizing search results, and organizing ideas. Palani et al. [54] use log analysis and self-reports in a study with 34 design students. They observed three main goals: To get an overview of the information space, to discover design patterns and criteria, and to seek inspiration and generate ideas. In the study they specifically took note of participants first struggling to find adequate terms to describe their information need, using rapid querying and query reformulating to converge on suitable terms. Furthermore, they observe that participants usually go through a divergent exploration phase before a convergent synthesis phase. Based on a previous online survey and study [80, 81], Li et al. [45] investigate information behavior in a diary study of 11 university members on self-chosen creative tasks. They employ Sawyer’s eight-stage creativity framework and specifically focus on the usage of resources (search, images, Q&A, social sites, videos). They grouped the observed information resource usage into five categories: look for specific information, support creative processes, learn domain knowledge, learn procedural knowledge, and manage (organize) found information. Specifically for images they distinguish more fine-granular uses (e.g., as per Pinterest, Instagram, Tumblr, Flickr, and image search): support ideation and other creative processes, see finished examples, find out what they themselves liked or disliked, and manage found information in an overview. They notice that image search engines were mostly used to look for a broad range of images, whereas image sites like Pinterest and Instagram were often used to browse for higher-quality images created by specific artists or professionals. To summarize, we find that, though different studies use different categories, there are three common themes in search for creative tasks: searching to learn, to get inspired, and to get an overview. We also see these themes in the use of text-to-image generation models, specifically in our case study (Section 4).

2.4 Interactive Retrieval

Interactive information retrieval focuses on users’ search or information-seeking behavior and, subsequently, developing new interaction methods that assist users in this process [60]. We review work on query log analysis (as a source to study the interactive behavior of users), and query reformulation and suggestion as two frequent interaction methods.

Query Log Analysis. The analysis of query logs was already a valuable resource in the early days of web search, e.g., to improve the retrieval effectiveness [36], or to study the behavior of users [7, 32–34]. For instance, Broder [7]

¹The Midjourney community on the Discord platform counts more than 3 million members as of October 2022.

established the taxonomy of query intents into informational, navigational, and transactional information needs, and this taxonomy is still used today [1]. Query logs are also valuable for studying the interactive behavior of users. For instance, Jansen et al. studied the query reformulations conducted by searchers, finding that reformulations are either generalizations (subset of words), specialization (superset of words), synonymous, or on a different topic.

Query Reformulation. Query reformulation approaches aim to improve retrieval effectiveness by replacing the original query with substituted or expanded reformulations [15]. Thereby, query reformulation can be either precision-oriented (when a term is substituted by a more specific one) or recall-oriented (if the query is expanded). Searchers do not start with perfect queries and instead reformulate: more than 50% of searchers reformulate at least one query during a search session [34]. Automatic query expansion approaches, like RM3 [31], may leverage (pseudo) relevance feedback to add new (weighted) terms to the initial query and thereby address the vocabulary mismatch problem prominent in textual retrieval. However, it is not yet clear what reformulations are helpful in which situations during the creative work with text-to-image generation models (i.e., precision-oriented or recall-oriented reformulations).

Query Suggestion. Search engines assist their users and provide a list of suggested queries for some input query [5], which is referred to as query auto completion [8] if the query is not finished. Thereby, query suggestions are important, as 30% of the queries among a commercial query log were previously suggested to its users [17]. Spelling corrections are particularly important query suggestions, as 10–15% of search queries have spelling errors [14]. Furthermore, query suggestions often aim at supporting users by showing related terms [28], and query log analysis shows that those suggested related terms are heavily used [33]. However, it is important not to overwhelm users, as showing fewer alternatives for the suggestions is better than many [75]. Overall, our survey highlights that users highly value the interaction methods employed in “traditional” search or information-seeking situations. Consequently, framing text-to-image generation models in a search setting can enable similar benefits to users with creative tasks.

3 CASTING TEXT-TO-IMAGE GENERATION AS SEARCH

Users of text-to-image models do not know which prompts to use to steer resulting images in a particular direction, putting them in a situation resembling image search. Thereby, prompt engineering resembles query refinement in image search to “find” relevant images. However, instead of retrieving images for a query from a finite set of indexed images, a result is returned for every conceivable prompt as a query. Moreover, due to the possibility of generating random variations for images, a generative model is effectively an “infinite index.” Consequently, we conceptualize generative text-to-image models as an index integrated into a retrieval system: from the user perspective, the interaction methods are presented (Section 3.1), and from the system perspective, retrieval technologies are examined (Section 3.2). Finally, we present requirements for the evaluation of such systems (Section 3.3).

3.1 Interaction Methods

Though the characteristic way of interacting with text-to-image generation models is the text prompt, in a very short time other methods have been added to the respective interfaces to aid in the image generation process. To illustrate possibilities, we here capture a brief snapshot of interaction methods based on the most popular models at in October 2022. This collection does not include related models, like text-to-video or text-to-3D generation models [25, 55].

Prompting. When using text-to-image generation models, the characteristic interaction method is prompting the model. This interaction method serves as the first contact point with the model in the image generation process,

much like querying for regular web search. Some interfaces allow to include images in the prompt to further push the generated images into a certain direction. However, unlike in pure content-based image retrieval, the model interfaces usually require the prompt to contain some text. Another aspect of prompting in some interfaces is to specify model parameters along with the prompt, for example the size of the image to-be-generated or whether to generate tiling images. Like in web search, a negation operator can be used to exclude certain concepts from the to-be-generated image. The widespread open-source model Stable Diffusion only provides a command line interface, but its community has implemented several graphical interfaces for it, for example one maintained by AUTOMATIC1111² (cf. Figure 1 (a)).

Moreover, several services appeared in the larger text-to-image generation ecosystem to aid the user in prompt engineering. Specialized search engines allow to search for images that users generated with text-to-image generation models. The engines then reveal the prompts to generate the found images, thereby allowing for prompt-reuse. Images are indexed either by their prompt or by image content (e.g., using CLIP). Examples of such search engines include the “community feed” of the Midjourney web app or the independent search engine Lexica that indexes images from the Stable Diffusion Discord channel (cf. Figure 1 (b)). For the latter, according to its developer it served 1.4 million queries in a week and its index contains 12 million images in September 2022, and has raised USD 5 million, clearly indicating the need for such systems. Other services allow for (social) prompt engineering in a click interface³ or even to buy prompts that are said to produce consistent results.⁴ Yet other projects attempt careful analyses of how the prompt affects the result, providing huge lists of examples.⁵

Variations. When an image is generated, this interaction method allows to change parts of the composition of an image. This is useful if an image result is mostly satisfying, but need improvements in certain aspects. We can distinguish three types for variation generation: (1) the user does not change the prompt, causing only minor and random changes in the composition (cf. Figure 1 (d)); (2) the user changes the prompt, giving the model a new target as it continues from a generation checkpoint of the original image; (3) the user specifies a semantic edit of the image, changing elements of the original image while preserving its original characteristics [37].

In- and Out-Painting. When an image is generated, this interaction method allows to restrict variations to user-specified regions of the image. If the region is inside the original image, this method is called in-painting (cf. Figure 1 (c)), if the region is outside, it is called out-painting, but both work similarly: the model attempts to fill the region in a way that fits to the prompt and also to the parts of the original image at the border of the region. This kind of variations is useful if the user wants to change a specific region of the generated image (in-painting), or if the user wants to extend an image (out-painting).

Quality Enhancements. When the user is satisfied by the composition of an image, this interaction method allows to enhance the image quality in one or multiple ways without changing the composition. The most common type of quality enhancement is upscaling the image to a higher resolution. However, often different upscaling methods exist, which can create from the same source image other images that appear sharp or soft, realistic or artistic, all while keeping the original composition. Choosing a distinct-looking upscaling algorithm can be useful when creating several images that need to share a common feel. A different kind of enhancement is the use of image-to-image models

²<https://github.com/AUTOMATIC1111/stable-diffusion-webui>

³E.g., <https://phraser.tech>

⁴<https://promptbase.com>

⁵E.g., <https://github.com/willwulfken/MidJourney-Styles-and-Keywords-Reference>

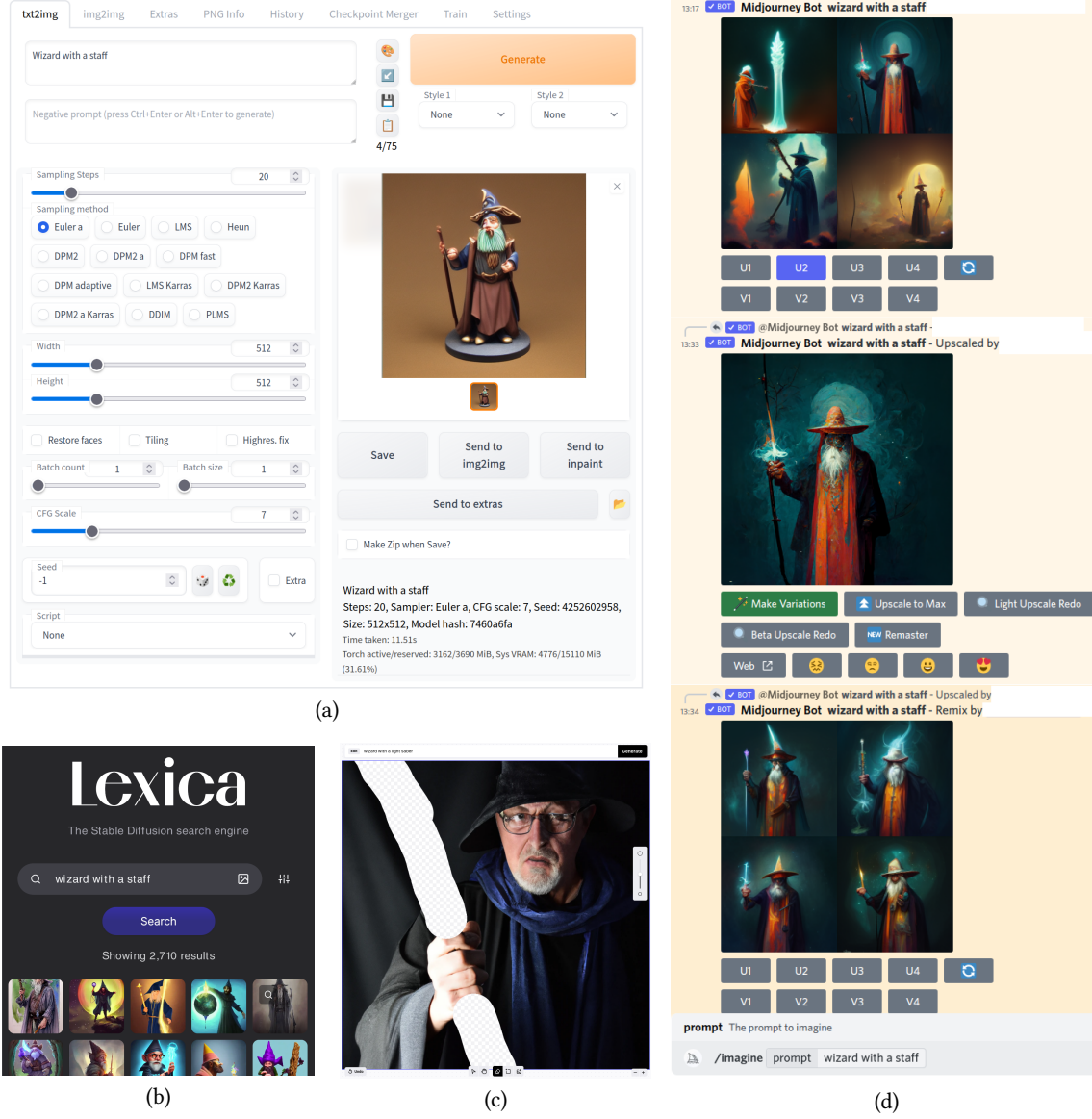


Fig. 1. Screenshots exemplifying interfaces and the interaction methods discussed in Section 3.1: (a) prompting in a community-maintained Stable Diffusion web interface; (b) Lexica search engine for generated images; (c) in-painting in DALL-E 2 on an image originally generated for the prompt “wizard with a staff”: the staff is manually masked and the prompt changed to “wizard with a light saber”; (d) upscaling and variation generation in Midjourney (anonymized).

specifically trained on face correction [74]. We expect that also other image-to-image models specialized on specific operations will be integrated in the future.

Image-to-Text. When the user wants to reformulate the prompt, but also use parts of generated images, this interaction method allows them to get a textual description of the image that reads like a prompt.

3.2 Employing Information Retrieval Technologies in Text-to-Image Generation

As the above list of interaction methods demonstrates, the text-to-image generation community created support for a multitude of interaction methods, but so far restricted information retrieval concepts to external tools. This section highlights how different retrieval technologies could be employed to support users in the generation of images.

A central aspect of information retrieval research is the concept of a result’s relevance. Like in regular image retrieval, the concept of relevance depends on the users’ information needs, of which seven different categories have been distilled from the literature (cf. Section 2). Generating images instead of finding some can, at least in theory, fulfil most of these needs, and may be particularly useful for needs of entertainment, illustrations, aesthetic appreciation, engaging others, inspiration, and social interactions. For social interactions, for example, it is very useful that generative models account for general moods mentioned in the prompt, providing a clear path towards generating images that evoke certain emotions. The remaining information need, for which generating images is in some cases not the right choice, is the need for knowledge construction, as generated images often do not account for real-world knowledge. In this regard, it is necessary to distinguish two different intents when generating images. First, the user may already have a clear vision of the target image, for example for an illustration. A user with this intent iteratively refines their prompt until the system generates an image close to their vision, which we refer to as descriptive approach. Second, they may not have a clear vision or not a clear goal, but only a set of constraints. Also, with this intent the user iteratively refines their prompt as feedback to random elements introduced by the system, but directing the system loosely towards an image of their liking, which we refer to as creative approach. Though very different from the user’s point of view, both approaches seem identical for the system: A generic prompt is expanded with details to become more specific.

Moreover, research of interactive retrieval (cf. Section 2) is also applicable for text-to-image retrieval. Query log analysis can be essential to identify keywords in prompts that yield satisfactory results in general, and to identify search missions and early abortions that can point to problems in the model. In these cases, we expect that methods of query suggestion can be very helpful to assist especially novice users. At the moment, reformulating prompts is challenging as users are not able to predict the effect of prompt changes on the generated images. In our case study (Section 4), the creative professional thus abstained from optimizing the prompt, rather trying completely new ones instead. We see here a clear lack of user support in retrieval terms in current interfaces. External tools like prompt search engines try to compensate for this deficiency, but will not be able to reach the effectiveness of integrated solutions as they are now widespread in search engines (see Section 5 for a discussion of possible remedies).

With these considerations in mind, one is able to adapt the notion of relevance and retrieval methods like query suggestion from image retrieval to text-to-image generation, and thus to employ IR evaluation measures.

3.3 IR Evaluation

Framing text-to-image generation as a search problem enables measuring the effectiveness of rankings produced by the end-to-end system following the standard experimental practices in IR. However, we show that the virtually never-ending “infinite index” has severe implications for the design and evaluation of experiments because the full recall base is unknown. Furthermore, we discuss the challenges that emerge from this unknown recall base for constructing re-usable benchmark collections and speculate on approaches to address those challenges. We focus on ranking effectiveness measures as other aspects, such as the interface design and layout, are ignored in Cranfield-style evaluations.

Implications of the Infinite Index on IR Evaluations. The virtually never-ending stream of images results in problems for recall-oriented evaluation measures and potentially overestimated effectiveness scores due to (near-) duplicated images. Effectiveness measures can be categorized into utility-oriented (based only on the ranking) and recall-oriented (normalized by some “best possible” ranking) evaluation measures [49]. For a virtually infinite number of images, it is unknown how many (highly) relevant images exist, so for recall-oriented measures like nDCG [35], the normalization is only straightforward in cases where already enough highly relevant images are identified so that the best possible ranking has highly relevant images on all positions. Utility-oriented measures (like Precision@k, MRR, RBP [51], etc.) do not suffer from this problem as they only measure the effectiveness of a ranking by the images available in the ranking. Given the virtually infinite number of images, retrieval approaches might rank many/only (near-) duplicated images that might be relevant in isolation but are non-distinguishable for users. Evaluation measures that operate on rankings with (near-)duplicates overestimate the effectiveness [4, 19], and also learning-to-rank approaches learn suboptimal ranking models when trained on redundant data [18]. Consequently, it is important to deduplicate rankings before evaluations accordingly, and utility-oriented measures (like RBP) on deduplicated rankings with complete judgments for the top-k images allow theoretically well-grounded evaluations even on the infinite index of text-to-image models.

Evaluations with Active Judgment Rounds. The experimental evaluation of retrieval systems is usually conducted according to the Cranfield paradigm [12, 13] that assumes that all documents are judged for all information needs. The original Cranfield experiments [12, 13] were conducted on a collection of 1,400 documents and complete relevance judgments for 225 topics. However, complete judgments became infeasible almost immediately thereafter since collection sizes grew substantially. The current best practice at shared tasks in IR is to create per-topic pools of the submitted systems’ top-ranked documents and then judge each topic’s pool [72], assuming that unjudged documents are non-relevant. However, this assumption that the judgment pools are “essentially complete” might be wrong in the case of an infinite index, especially if query expansion approaches are involved. Consequently, rigorous evaluations must involve costly judgment rounds of unjudged images to assume completeness (e.g., for the top-k results), at least for utility-oriented measures, which hinders fully automated evaluations.

Evaluations without Active Judgment Rounds. Research in IR largely benefited from the availability of robust and reusable test collections constructed during shared tasks [71]. However, those collections are only robust if most of the unjudged documents are non-relevant [72], which might not be true for text-to-image models. Consequently, the construction of robust and reusable test collections is a major challenge that would require the experience from multiple diverse shared tasks and subsequent post-hoc experimentations (e.g., some robustness checks for traditional test collections are conducted years after their construction [73]). Hence, all post-hoc experiments on an infinite index would require proper handling of unjudged images. Traditionally, unjudged documents are either simply removed (condensing a system’s result lists to the contained judged documents in their relative order) [62], assumed to be non-relevant (the default) or highly relevant (lower/upper bounds) [49], or their relevance label might be predicted [2]. While those approaches are well-studied for traditional retrieval experiments (e.g., condensed lists often overestimate effectiveness scores [63] and the gap between lower and upper bounds can be very large [49]), it is not yet clear if they are appropriate for an infinite index. Consequently, it is not yet clear how robust and reusable test collections have to be constructed, but we speculate that techniques from machine translation (e.g., measuring the similarity of an unjudged document via phash [77] to judged reference images) or relevance prediction might be suitable.

4 CASE STUDY OF TEXT-TO-IMAGE GENERATION FOR CREATIVE TASKS: GAME ARTWORK SEARCH

To illustrate information retrieval on generative text-to-image models, we report on a case study in which such a model is employed in a creative task. Section 4.1 describes the study setup and exemplary creative task, the generation of artwork for an online card game. Section 4.2 then summarizes the main observations from the study. A full report of the study is available as supplementary material.⁶

4.1 Setup of the Case Study

For the case study we recruited via personal contacts one creative professional who allowed us to observe him as he explored the usage of text-to-image generation models in his creative process. The professional described himself as generalist game designer and developer with an experience of five major game releases, and as teacher for game development at a university. Before the case study, he described himself as being very intrigued by text-to-image generation models, which he encountered in his Twitter feed, having also already watched some online videos on the technology (“2 minute papers”). He had already generated about 50 images in DALL-E 2, about 20 in Midjourney, and less than 10 using Stable Diffusion on his own hardware, but none of these as part of a project. But he assumed that text-to-image generation models would be very useful for the video game industry.⁷

Based on his experience, the professional chose to explore for the study the use of text-to-image generation models in artwork creation for an online card game. Specifically, he was interested in creating a “deck-building online card game like Magic the Gathering set in a fantasy universe.” In the game, each playing card has its own artwork that visually links it to the fantasy universe. Furthermore, the cards belong to different “factions,” which need to be visually distinguishable. The professional decided up front to work towards a “concept art like style.” In the 5 hours we allocated for the study, the professional expected to first create a “mood board” of images to capture the artistic style of the desired artwork [42] and then to create the artwork itself for a few cards. Based on his own tests he decided to employ Midjourney for this task. This choice reflects the concept of Midjourney, which emphasizes “painterly aesthetics” and aims to assist creators especially at project start to “converge on the idea they want much more quickly” [64].

The case study was conducted using the think-aloud method, with additional questions being asked as the professional was waiting for images to be generated. Since the study did not target interface design, one of the authors extensively used Midjourney in preparation for the study to provide technical assistance for the professional. To persist observations we employed extensive note-taking, video and audio recording, and the logging facilities of the Midjourney web app. Following Li et al. [45] we employed forms to structure our notes for different events, in our case for queries, troubles, and shifts in design goals. A report on the study including all generated images is available as supplementary material.

4.2 Main Observations from the Case Study

This section summarizes the insights from the case study as three main observations. Specifically, we found the mood board to be a central tool for the professional, and analyze its use as per the five reasons for using information resources [45]. For an analysis of the professional’s mental state we employ Kuhlthau’s model of the information search process [40]. And based on the professional’s comments during the study we identified his mentioned lack of control as major problem that needs to be addressed by future tools.

⁶<https://doi.org/10.5281/zenodo.7221435> (anonymized for blind review)

⁷Video games account for about 57% of revenue in the Digital Media market in 2022, namely US\$197 billion [29]

The mood board as prompt library. Lemarchand [42] defines a mood board as “a single page or screen of pictures arranged around a certain idea or theme” that serves two main purposes: first to inspire new ideas sparked by placing images next to each other (support creative processes), and second to convey a concept quickly and effectively (manage found information). Having created the mood board out of images from the community feed of Midjourney, however, the professional directly started to use the mood board also as source for his prompts. When creating a new image, he selected from the mood board the image closest in art style to what he had in mind and copied the “style part” of that image’s prompt for his own creation (cf. gray text in Figure 2). He thus used the mood board additionally to learn domain knowledge (style names, render engines, ...) and procedural knowledge (parameters like “--q 2” to increase the image quality). The professional only once used an external search engine to look for the name of artists that worked on “Magic the Gathering” cards, and was pleased to find out that these were already part of the prompts he had copied. However, learning happened on a very shallow level, with the whole style-part of the prompt being copied and used like an atomic unit. This way of treating the style-part as an unit seems to be common in the text-to-image generation community at the moment, with even commercial services appearing which sell such units.⁸

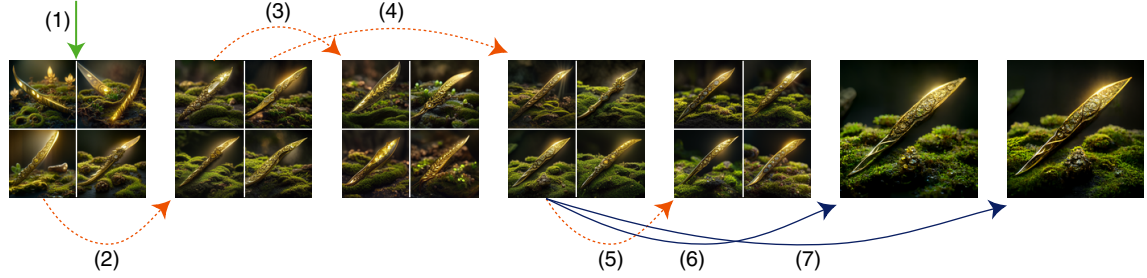
Uncertainty never fully ceases. In Kuhlthau’s model of the information search process, the searcher moves from uncertainty to understanding as the search progresses [40]. We observed clear parallels to this model and its stages during the case study, especially the stages of selection, exploration, formulation, and collection. In the selection stage the professional uses the mood board as inspiration to select content and style for a new image. In the exploration stage he created and modified the prompt: he mentioned strong feelings of uncertainty about what results he would get and how to modify the prompt to get what he had in mind. Having found something he deemed promising, he entered the formulation stage and focussed on repeatedly generating variations, becoming clear on certain aspects the final image should have. With a clear sense of direction, he would in the collection stage then upscale fitting images, testing the different upscale algorithms if necessary. As accounted for in the model, the professional went also back to earlier stages, specifically if he saw a dead end (cf. Figure 2). However, Kuhlthau mentions two “types of uncertainty,” and though uncertainty regarding the concept (what he is looking for) diminishes as described above, uncertainty regarding the technical process (how to get there) stays high with the AI staying largely unpredictable for him.

Sense of direction, but lack of control. Though the professional remarks in some situations that the inherent unpredictability of the process is appealing (“I also wanted to be surprised”), he also mentions that the process is very strenuous, which we link to him often going back in the history of his generated images to assess over and over again which interactions produced good results and from which image to continue. An interface that supports the user in organizing generated images thus seems necessary. Though the professional remarked that he developed a good intuition of the direction into which image variations develop, he also felt a lack of control. He said he decided whether to continue down a path or try another one, but did not feel like he could change the direction. After the case study, Midjourney introduced the possibility to modify the prompt as variations are generated, but the professional tried it and remarked that this does not solve the problem of selecting the right words for the prompt. Therefore, we find that the uncertainty of how to change the prompt to get the desired results has a large negative impact on the user’s sense of control.

The case study did indeed show clear parallels between text-to-image generation and image search. Specifically, we found existing theoretical models of the (creative) search process to be mostly applicable. The main difference we find

⁸E.g., <https://promptbase.com>

Initial prompt: an ancient golden dagger lying on moss, illuminated by godrays, close up, digital painting, matte painting, midjourney, concept art, detailed art, sci-fi cinematic painting, magic the Gathering, volumetric light, masterpiece, volumetric realistic render, epic scene, 8k, post-production detailed art, sci-fi cinematic painting --q 2



Reformulated prompt: a medieval dagger lying on moss, lit by god rays, art by Adrian Smith + Paul bonner, magic gathering style, warcraft, blizzard style, hearthstone, fantasy concept art, medieval, masterpiece, mystical, witchcraft

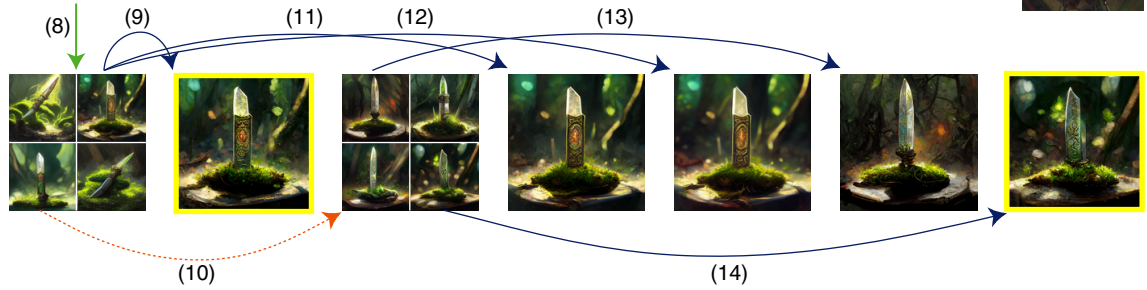


Fig. 2. Example image generation process from the case study, consisting of 14 steps taken in 22 minutes. Gray prompt text is copied from the prompt of another image in the mood board. For the reformulated prompt—after abandoning the first series of images as “leading nowhere”—, it is copied from the image that is part of the prompt. Interactions are text-to-image generation of four images (\rightarrow), generating four variations of one image (same prompt, $\cdots\rightarrow$), and upscaling one image (\longrightarrow). The “beta” upscaling method is used in step 12, the “light” method in step 13, and the default (“detailed”) in all other cases. The professional selected the two images with yellow border at the end. Even though he remarked the image generated in step 9 does not show a dagger as he intended, he found it intriguing, saying it evoked a story in his head, especially in combination with the other selected image.

is located in the not-ceasing uncertainty on how to get to a certain result—though, the index being virtually infinite, the user has to assume that there is a way to get there. Based on our observations, we believe that tools that provide the user with more intuitive means to control the generation process are needed and can close this gap.

5 DISCUSSION

In this section, we discuss the limitations of generative text-to-image models, address how active learning could help to guide the generation and identify ethical concerns.

5.1 Limitations

While the functionality provided by text-to-image models is already of sufficient quality to be used in real-world applications [52, 64], we identified multiple limitations related to the workflow or the possibilities given by current methods—the same workflow that was also used in the case study.

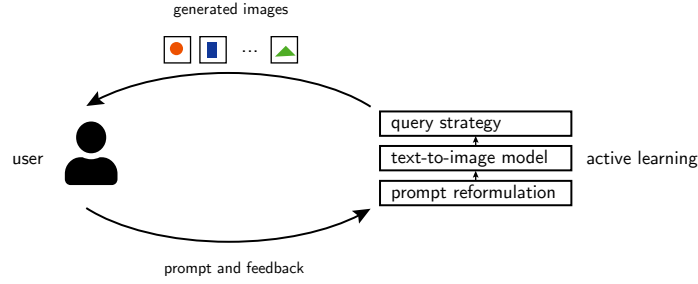


Fig. 3. A conceptual overview of the active learning loop for the use case of guided text-to-image generation.

Prompt Engineering. Although prompt engineering has been the method of choice in a similar recent setup where screenplays and theatre scripts are co-written by a writer and a large language model [50], it is not a good interface for text-to-image scenarios. Users have quickly realized that iteratively adding modifiers to the prompt (as done in Section 4), causing the model to apply desired outcome styles to the generated image, is the most effective way to guide the image creation process [47, 53]. This has given rise to a whole new subfield of text-to-image prompt engineering [47] in which prompts have become rather long chains of keywords instead of textual descriptions, strongly resembling search engine optimization where by carefully selecting and stuffing keywords the users have learned to adapt to the algorithm instead of vice versa.

Influence of the Training Data. One fundamental limitation of the current models is that both text encoders and diffusion models generate new data by blending together concepts that have been learned from large datasets. If we write a prompt containing a concept that does never or rarely occur either in the text corpora or image datasets, it will likely to lead to sub-par generated image. One could of course argue that unknown concepts can be described as a paraphrase, e.g. assuming the training data had not contained a single image of a centaur, a fitting prompt describing a “mythical creature with a horse’s body and a man’s torso” might lead to the desired outcome rather fast. The extent of this limitation must be investigated in future research.

5.2 Active Learning to Guide Text-to-Image Generation via Relevance Feedback

From the case study we have learned that targeted text-to-image generation is already surprisingly effective and can yield highly remarkable results. As presented in Section 3.1 the current manner of operation boils down to iterative prompt engineering which in turn is a fundamental limitation as outlined in Section 5.1. In this section we propose active learning as a solution to this problem and outline how it can be integrated as a feedback mechanism into a Stable Diffusion workflow.

Active learning [43, 79] is an iterative procedure between a user and a machine learning model which is used in order to train a model from user input when no training data is available. The goal is to maximize the model quality while minimizing the required user input. A full active learning setup consist of a model *that is trained on some task*, a query strategy *that selects data from an existing resource or generates new data to be labeled*, and a stopping criterion *that indicates when the process is likely to stop improving the result* [65]. During each iteration the query strategy provides examples to the user who annotates them according to the respective task. After each iteration, a new model is trained on all data labeled so-far and the loop is repeated until the user is satisfied with one of the generated images.

Text-to-Image Generation with Relevance Feedback. For text-to-image generation, the overall active learning setup is shown in Figure 3. The process starts with the user and an initial prompt. The active learning model learns to reformulate prompts which are in turn passed to the text-to-image model. The model is trained with the user feedback as target values so that the resulting images should increasingly become more appealing to the user. Next, the query strategy decides which images are shown to the user. It balances exploration versus exploitation, a tradeoff that is well-known in information retrieval [3, 26] where exploration means selecting images that are different from the current best candidates, and exploitation means selecting images that are close to the current best solutions. Finally, the stopping criterion is the user who stops the process once his information need is satisfied. We use active learning as a relevance feedback mechanism as has been successfully done in numerous previous work [69, 76, 79].

Feedback Types. There are different types of user feedback that can be applied to obtain relevance feedback: (1) binary relevance feedback [22] where the user rates each image in relation to the target concept as “unappealing” or “appealing”; (2) graded relevance feedback [22] where the user grades each image on a multi-point scale (e.g., from 0 to 5) from unappealing to appealing; (3) ranking where the user ranks each image (possibly also images from previous iterations) from unappealing to appealing. Moreover, and specific to text-to-image generation, the user can update the prompt during each iteration.

Challenges. The key challenge for this feedback mechanism is to revert the images back to a textual representation, which then can be used to learn the prompt reformulation. We could of course use latent image representation as input for another model to predict the user’s response, however, this would tackle the problem completely in the image space, and thereby we would mostly ignore the textual embedding space. Getting a textual description from a given image, i.e. image-to-text models, have recently gained increased attention [20, 21] but still have some drawbacks such as requiring multiple images to create one text [20]. Once this reverse direction is improved, this will open up the full range of natural language processing and information retrieval methods to effectively process user feedback to improve prompts during the reformulation step.

When text-to-image generation is cast as a search problem (as shown in Section 3), the process of trying different prompts until a satisfying image is generated becomes similar to traditional image retrieval and therefore incorporating active learning as a relevance feedback mechanism is an obvious choice of a well-established method. We predict that once image-to-text will reach sufficient maturity, active prompt generation will become the favored interface between a user and a text-to-image generation model (disregarding editing capabilities such as in-painting or out-painting which are orthogonal to this approach).

5.3 Ethical Concerns

A paradigm powerful enough to generate artifacts like images, texts, and other media types that are sometimes difficult to distinguish from human illustrations raises ethical concerns.

Will algorithms replace artists? We begin with the obvious question: will text-to-image generation models threaten the jobs of artists? First, drawing from the experiences we gained through the case study, it is currently difficult to guide text-to-image models towards a desired outcome. Moreover, the decision whether the generated images depicts the desired scene and is of sufficient quality still needs to be made by the user. Therefore, we think this new paradigm will be a powerful tool but not replace the human illustrator for the foreseeable future—even assuming the image quality might eventually reach human levels, which is clearly not the case yet. This is corroborated by others such

as Liu et al. [48], who developed and evaluated a system that supports the user at generating illustrations for news articles, finding that artistic knowledge is still beneficial to the generated result and explicitly stating that “generative AI deployment should [...] augment rather than [...] replace human creative expertise”. We support this view: Instead of an autonomous AI that acts of its own accord, a “supportive AI” that asks and incorporates the decisions made their users. In this regard, a retrieval system by definition, is a supportive AI-based system, since they traditionally act only on behalf of users who have an information need.

Who is the author of a generated image? And who owns the rights? This is a currently unresolved situation causing uncertainty regarding the use of AI-generated images, which is why large platforms such as the well-known stock image provider Getty Images have recently banned all AI-generated content.⁹ In the end, this decision must be made by policy makers, if not by courts of law where past copyright disputes set many legal precedents.

Can text-to-image models be abused for automated or non-automated misinformation? Generated misinformation is already a pervasive problem and widely discussed in the context of so-called “deep fakes” and AI-generated text generation [39, 67, 78]. To mitigate this problem in text-to-image models like Stable Diffusion, image watermarking is employed, marking an image as being artificially generated.¹⁰ Though watermarks are not easily removed, this may not be enough, unless they are checked on virtually all devices. This, however, will require policy makers to legally require device makers to check and warn users. Moreover, watermarking images raises privacy concerns of its own.

Do these models express or even amplify bias? Bias in language models is a known problem [38], too. Similarly, text-to-image models will be biased, since their training data distribution has not been systematically checked for social and other kinds of biases. In this regard, retrieval process on top of generative models may also be used to mitigate their inherent biases. For example, in information retrieval research on web search, fair ranking is meanwhile a widely studied problem. Image search engines based on generative models will have to postprocess and re-rank their results to control for biases, just like their traditional counterparts. However, the technologies applied for traditional search engines can be directly transferred to search engines based on generative models.¹¹

All in all, the proposed approach might reinforce existing ethical issues that are known for text encoder models and image generation models but does not introduce any additional problems. Nevertheless, one has to be aware of these issues and they have to be kept in mind during future research.

6 CONCLUSION

The use of generative text-to-image models will require facilitating systems and services. Their integration into existing systems is already ongoing, as can be observed for years already with writing assistance and translation systems, but now also in more creative areas. An integration into end user software to create slide presentation or artwork, however, will not fully serve the needs for everyone who searches for inspirational imagery. In this respect, dedicated search engines based on generative text-to-image models as indexes, with a dedicated user interface to formulate information needs, and tailored retrieval models, are likely already under development. However, the development of a search engine is not trivial, and the information retrieval community is facing a renewed challenge to develop an understanding and the technological basis for such search engines. This includes the development of new retrieval models and relevance

⁹<https://voicebot.ai/2022/09/23/getty-images-removes-and-bans-ai-generated-art/>

¹⁰<https://github.com/CompVis/stable-diffusion>

¹¹ Compare for example *lexica.art*’s result (<https://lexica.art/?q=nurse>) with that of Google Images (<https://www.google.com/search?tbm=isch&q=nurse>)

scores as well as adapting evaluation methodologies to benchmark search engines based on generative models. Finally, what applies to generative text-to-image models will likely also apply to generative models of other kinds, which opens a whole new world of exciting new research directions and promises high impact.

REFERENCES

- [1] Daria Alexander, Wojciech Kusa, and Arjen P. de Vries. 2022. ORCAS-I: Queries Annotated with Intent using Weak Supervision. In *SIGIR '22: The 45th International ACM SIGIR Conference on Research and Development in Information Retrieval, Madrid, Spain, July 11 - 15, 2022*, Enrique Amigó, Pablo Castells, Julio Gonzalo, Ben Carterette, J. Shane Culpepper, and Gabriella Kazai (Eds.). ACM, 3057–3066. <https://doi.org/10.1145/3477495.3531737>
- [2] Javed A. Aslam and Emine Yilmaz. 2007. Inferring document relevance from incomplete information. In *Proceedings of the Sixteenth ACM Conference on Information and Knowledge Management, CIKM 2007, Lisbon, Portugal, November 6-10, 2007*, Mário J. Silva, Alberto H. F. Laender, Ricardo A. Baeza-Yates, Deborah L. McGuinness, Bjørn Olstad, Øystein Haug Olsen, and André O. Falcão (Eds.). ACM, 633–642.
- [3] Kumaripaba Athukorala, Alan Medlar, Antti Oulasvirta, Giulio Jacucci, and Dorota Glowacka. 2016. Beyond Relevance: Adapting Exploration/Exploitation in Information Retrieval. In *Proceedings of the 21st International Conference on Intelligent User Interfaces, IUI 2016, Sonoma, CA, USA, March 7-10, 2016*, Jeffrey Nichols, Jalal Mahmud, John O'Donovan, Cristina Conati, and Massimo Zancanaro (Eds.). ACM, 359–369. <https://doi.org/10.1145/2856767.2856786>
- [4] Yaniv Bernstein and Justin Zobel. 2005. Redundant documents and search effectiveness. In *Proceedings of the 2005 ACM CIKM International Conference on Information and Knowledge Management, Bremen, Germany, October 31 - November 5, 2005*, Otthein Herzog, Hans-Jörg Schek, Norbert Fuhr, Abdur Chowdhury, and Wilfried Teiken (Eds.). ACM, 736–743. <https://doi.org/10.1145/1099554.1099733>
- [5] Sumit Bhatia, Debapriyo Majumdar, and Prasenjit Mitra. 2011. Query suggestions in the absence of query logs. In *Proceeding of the 34th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR 2011, Beijing, China, July 25-29, 2011*, Wei-Ying Ma, Jian-Yun Nie, Ricardo Baeza-Yates, Tat-Seng Chua, and W. Bruce Croft (Eds.). ACM, 795–804. <https://doi.org/10.1145/2009916.2010023>
- [6] Andrew Brock, Jeff Donahue, and Karen Simonyan. 2018. Large scale GAN training for high fidelity natural image synthesis. *arXiv preprint arXiv:1809.11096* (2018).
- [7] Andrei Z. Broder. 2002. A taxonomy of web search. *SIGIR Forum* 36, 2 (2002), 3–10. <https://doi.org/10.1145/792550.792552>
- [8] Fei Cai and Maarten de Rijke. 2016. A Survey of Query Auto Completion in Information Retrieval. *Found. Trends Inf. Retr.* 10, 4 (2016), 273–363. <https://doi.org/10.1561/15000000055>
- [9] Catherine Chavula, Yujin Choi, and Soo Young Rieh. 2022. Understanding Creative Thinking Processes in Searching for New Ideas. In *ACM SIGIR Conference on Human Information Interaction and Retrieval* (Regensburg, Germany). ACM, New York, NY, USA, 321–326. <https://doi.org/10.1145/3498366.3505783>
- [10] Hyerin Cho, Minh TN Pham, Katherine N. Leonard, and Alex C. Urban. 2021. A systematic literature review on image information needs and behaviors. *Journal of Documentation* 78, 2 (2021), 207–227.
- [11] Youngok Choi. 2013. Analysis of image search queries on the web: Query modification patterns and semantic attributes. *Journal of the American Society for Information Science and Technology* 64, 7 (2013), 1423–1441.
- [12] Cyril Cleverdon. 1967. The Cranfield tests on index language devices. In *Aslib proceedings*. MCB UP Ltd. (Reprinted in Readings in Information Retrieval, Karen Sparck-Jones and Peter Willett, editors, Morgan Kaufmann, 1997), 173–192.
- [13] Cyril W. Cleverdon. 1991. The Significance of the Cranfield Tests on Index Languages. In *Proceedings of the 14th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, Chicago, Illinois, USA, October 13-16, 1991 (Special Issue of the SIGIR Forum)*, Abraham Bookstein, Yves Chiaramella, Gerard Salton, and Vijay V. Raghavan (Eds.). ACM, 3–12.
- [14] Silviu Cucerzan and Eric Brill. 2004. Spelling Correction as an Iterative Process that Exploits the Collective Knowledge of Web Users. In *Proceedings of the 2004 Conference on Empirical Methods in Natural Language Processing, EMNLP 2004, A meeting of SIGDAT, a Special Interest Group of the ACL, held in conjunction with ACL 2004, 25-26 July 2004, Barcelona, Spain*. ACL, 293–300. <https://aclanthology.org/W04-3238/>
- [15] Van Dang and W. Bruce Croft. 2010. Query reformulation using anchor text. In *Proceedings of the Third International Conference on Web Search and Web Data Mining, WSDM 2010, New York, NY, USA, February 4-6, 2010*, Brian D. Davison, Torsten Suel, Nick Craswell, and Bing Liu (Eds.). ACM, 41–50. <https://doi.org/10.1145/1718487.1718493>
- [16] Prafulla Dhariwal and Alexander Nichol. 2021. Diffusion models beat gans on image synthesis. *Advances in Neural Information Processing Systems* 34 (2021), 8780–8794.
- [17] Alan Feuer, Stefan Savev, and Javed A. Aslam. 2007. Evaluation of phrasal query suggestions. In *Proceedings of the Sixteenth ACM Conference on Information and Knowledge Management, CIKM 2007, Lisbon, Portugal, November 6-10, 2007*, Mário J. Silva, Alberto H. F. Laender, Ricardo A. Baeza-Yates, Deborah L. McGuinness, Bjørn Olstad, Øystein Haug Olsen, and André O. Falcão (Eds.). ACM, 841–848. <https://doi.org/10.1145/1321440.1321556>
- [18] Maik Fröbe, Janek Bevendorff, Jan Heinrich Reimer, Martin Potthast, and Matthias Hagen. 2020. Sampling Bias Due to Near-Duplicates in Learning to Rank. In *43rd International ACM Conference on Research and Development in Information Retrieval (SIGIR 2020)*. ACM, 1997–2000. <https://doi.org/10.1145/3397271.3401212>

- [19] Maik Fröbe, Jan Philipp Bittner, Martin Potthast, and Matthias Hagen. 2020. The Effect of Content-Equivalent Near-Duplicates on the Evaluation of Search Engines. In *Advances in Information Retrieval. 42nd European Conference on IR Research (ECIR 2020) (Lecture Notes in Computer Science, Vol. 12036)*, Joemon M. Jose, Emine Yilmaz, João Magalhães, Pablo Castells, Nicola Ferro, Mário J. Silva, and Flávio Martins (Eds.). Springer, Berlin Heidelberg New York, 12–19. https://doi.org/10.1007/978-3-030-45442-5_2
- [20] Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit H. Bermano, Gal Chechik, and Daniel Cohen-Or. 2022. An Image is Worth One Word: Personalizing Text-to-Image Generation using Textual Inversion. *CoRR* abs/2208.01618 (2022). <https://doi.org/10.48550/arXiv.2208.01618> arXiv:2208.01618
- [21] Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit H. Bermano, Gal Chechik, and Daniel Cohen-Or. 2022. Imagen Video: High Definition Video Generation with Diffusion Models. *CoRR* abs/2208.01618 (2022). <https://doi.org/10.48550/arXiv.2208.01618> arXiv:2208.01618
- [22] Gregory Gay, Sonia Haiduc, Andrian Marcus, and Tim Menzies. 2009. On the use of relevance feedback in IR-based concept location. In *25th IEEE International Conference on Software Maintenance (ICSM'09)*. IEEE Computer Society, 351–360. <https://doi.org/10.1109/ICSM.2009.5306315>
- [23] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. 2020. Generative adversarial networks. *Commun. ACM* 63, 11 (2020), 139–144.
- [24] William S Hemmig. 2008. The information-seeking behavior of visual artists: a literature review. *Journal of documentation* (2008).
- [25] Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey A. Gritsenko, Diederik P. Kingma, Ben Poole, Mohammad Norouzi, David J. Fleet, and Tim Salimans. 2022. Imagen Video: High Definition Video Generation with Diffusion Models. *CoRR* abs/2210.02303 (2022). <https://doi.org/10.48550/arXiv.2210.02303> arXiv:2210.02303
- [26] Katja Hofmann, Shimon Whiteson, and Maarten de Rijke. 2013. Balancing exploration and exploitation in listwise and pairwise online learning to rank for information retrieval. *Inf. Retr.* 16, 1 (2013), 63–90. <https://doi.org/10.1007/s10791-012-9197-9>
- [27] Vera Hollink, Theodora Tsikrika, and Arjen P. de Vries. 2011. Semantic search log analysis: A method and a study on professional image search. *J. Assoc. Inf. Sci. Technol.* 62 (2011), 691–713.
- [28] Chien-Kang Huang, Lee-Feng Chien, and Yen-Jen Oyang. 2003. Relevant term suggestion in interactive web search based on contextual information in query session logs. *J. Assoc. Inf. Sci. Technol.* 54, 7 (2003), 638–649. <https://doi.org/10.1002/asi.10256>
- [29] Statista Inc. 2022. Digital Media Report - Video Games. <https://www.statista.com/study/39310/video-games/>.
- [30] Sethurathienam Iyer, Shubham Chaturvedi, and Tirtharaj Dash. 2017. Image Captioning-Based Image Search Engine: An Alternative to Retrieval by Metadata. In *Soft Computing for Problem Solving (SocProS'17) (Advances in Intelligent Systems and Computing, Vol. 817)*, Jagdish Chand Bansal, Kedar Nath Das, Atulya Nagar, Kusum Deep, and Akshay Kumar Ojha (Eds.). Springer, 181–191. https://doi.org/10.1007/978-981-13-1595-4_14
- [31] Nasreen Abdul Jaleel, James Allan, W. Bruce Croft, Fernando Diaz, Leah S. Larkey, Xiaoyan Li, Mark D. Smucker, and Courtney Wade. 2004. UMass at TREC 2004: Novelty and HARD. In *Proceedings of the Thirteenth Text REtrieval Conference, TREC 2004, Gaithersburg, Maryland, USA, November 16-19, 2004 (NIST Special Publication, Vol. 500-261)*, Ellen M. Voorhees and Lori P. Buckland (Eds.). National Institute of Standards and Technology (NIST). <http://trec.nist.gov/pubs/trec13/papers/umass.novelty.hard.pdf>
- [32] Bernard Jansen, D. Booth, and A. Spink. 2009. Patterns of Query Reformulation During Web Searching. *J. Assoc. Inf. Sci. Technol.* 60, 7 (2009), 1358–1371. <https://doi.org/10.1002/asi.21071>
- [33] Bernard Jansen, Amanda Spink, and Sherry Koshman. 2007. Web searcher interaction with the Dogpile.com metasearch engine. *J. Assoc. Inf. Sci. Technol.* 58, 5 (2007), 744–755. <https://doi.org/10.1002/asi.20555>
- [34] Bernard Jansen, Amanda Spink, and Jan Pedersen. 2005. A temporal comparison of AltaVista Web searching. *J. Assoc. Inf. Sci. Technol.* 56, 6 (2005), 559–570. <https://doi.org/10.1002/asi.20145>
- [35] Kalervo Järvelin and Jaana Kekäläinen. 2002. Cumulated gain-based evaluation of IR techniques. *ACM Trans. Inf. Syst.* 20, 4 (2002), 422–446.
- [36] Thorsten Joachims and Filip Radlinski. 2007. Search Engines that Learn from Implicit Feedback. *Computer* 40, 8 (2007), 34–40. <https://doi.org/10.1109/MC.2007.289>
- [37] Bahjat Kavar, Shiran Zada, Oran Lang, Omer Tov, Huiwen Chang, Tali Dekel, Inbar Mosseri, and Michal Irani. 2012. Imagic: Text-Based Real Image Editing with Diffusion Models. *CoRR* abs/2210.09276 (2012). arXiv:2210.09276 <http://arxiv.org/abs/2210.09276>
- [38] Hannah Rose Kirk, Yennie Jun, Filippo Volpin, Haider Iqbal, Elias Benussi, Frederic Dreyer, Aleksandar Shtedritski, and Yuki Asano. 2021. Bias Out-of-the-Box: An Empirical Analysis of Intersectional Occupational Biases in Popular Generative Language Models. In *Advances in Neural Information Processing Systems*, M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan (Eds.), Vol. 34. Curran Associates, Inc., 2611–2624. <https://proceedings.neurips.cc/paper/2021/file/1531beb762df4029513ebf9295e0d34f-Paper.pdf>
- [39] Sarah Kreps, R. Miles McCain, and Miles Brundage. 2022. All the News That's Fit to Fabricate: AI-Generated Text as a Tool of Media Misinformation. *Journal of Experimental Political Science* 9, 1 (2022), 104–117. <https://doi.org/10.1017/XPS.2020.37>
- [40] Carol Collier Kuhlthau. 1993. A Principle of Uncertainty for Information seeking. *J. Documentation* 49, 4 (1993), 339–355. <https://doi.org/10.1108/eb026918>
- [41] Lo Lee, Melissa G. Ocepik, Stephann Makri, George Buchanan, and Dana McKay. 2019. Getting creative in everyday life: Investigating arts and crafts hobbyists' information behavior. *Proceedings of the Association for Information Science and Technology* 56, 1 (2019), 703–705. <https://doi.org/10.1002/pra2.141> arXiv:https://asistdl.onlinelibrary.wiley.com/doi/pdf/10.1002/pra2.141
- [42] Richard Lemarchand. 2021. *A Playful Production Process: For Game Designers (and Everyone)*. MIT Press, Cambridge, MA.
- [43] David D. Lewis and William A. Gale. 1994. A Sequential Algorithm for Training Text Classifiers. In *Proceedings of the 17th Annual International ACM-SIGIR Conference on Research and Development in Information Retrieval*, W. Bruce Croft and C. J. van Rijsbergen (Eds.). Springer, ACM/Springer,

- 3–12. https://doi.org/10.1007/978-1-4471-2099-5_1
- [44] Xiaoqing Li, Jiansheng Yang, and Jinwen Ma. 2021. Recent developments of content-based image retrieval (CBIR). *Neurocomputing* 452 (2021), 675–689.
- [45] Yuan Li, Yinglong Zhang, and Robert Capra. 2022. Analyzing Information Resources That Support the Creative Process. In *Proceedings of the 2022 ACM SIGIR Conference on Human Information Interaction and Retrieval (CHIIR'22)* (Regensburg, Germany). ACM, 180–190. <https://doi.org/10.1145/3498366.3505817>
- [46] Wen-Cheng Lin, Yih-Chen Chang, and Hsin-Hsi Chen. 2004. From Text to Image: Generating Visual Query for Image Retrieval. In *Multilingual Information Access for Text, Speech and Images, 5th Workshop of the Cross-Language Evaluation Forum, CLEF 2004, Bath, UK, September 15-17, 2004, Revised Selected Papers (Lecture Notes in Computer Science, Vol. 3491)*, Carol Peters, Paul D. Clough, Julio Gonzalo, Gareth J. F. Jones, Michael Kluck, and Bernardo Magnini (Eds.). Springer, 664–675. https://doi.org/10.1007/11519645_65
- [47] Vivian Liu and Lydia B Chilton. 2022. Design Guidelines for Prompt Engineering Text-to-Image Generative Models. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems* (New Orleans, LA, USA) (CHI '22). Association for Computing Machinery, New York, NY, USA, Article 384, 23 pages. <https://doi.org/10.1145/3491102.3501825>
- [48] Vivian Liu, Han Qiao, and Lydia Chilton. 2022. Opal: Multimodal Image Generation for News Illustration. *arXiv preprint arXiv:2204.09007* (2022).
- [49] Xiaolu Lu, Alistair Moffat, and J. Shane Culpepper. 2016. The effect of pooling and evaluation depth on IR metrics. *Inf. Retr. J.* 19, 4 (2016), 416–445.
- [50] Piotr Mirowski, Kory W. Mathewson, Jaylen Pittman, and Richard Evans. 2022. Co-Writing Screenplays and Theatre Scripts with Language Models: An Evaluation by Industry Professionals. *arXiv preprint arXiv:2209.14958* (2022).
- [51] Alistair Moffat and Justin Zobel. 2008. Rank-biased precision for measurement of retrieval effectiveness. *ACM Trans. Inf. Syst.* 27, 1 (2008), 2:1–2:27.
- [52] OpenAI. 2022. DALL-E: Creating Images from Text. <https://openai.com/blog/dall-e/>.
- [53] Jonas Oppenlaender. 2022. Prompt Engineering for Text-Based Generative Art. *arXiv preprint arXiv:2204.13988* (2022).
- [54] Srishti Palani, Zijian Ding, Stephen MacNeil, and Steven P. Dow. 2021. The "Active Search" Hypothesis: How Search Strategies Relate to Creative Learning. In *Proceedings of the 2021 Conference on Human Information Interaction and Retrieval* (Canberra ACT, Australia). ACM, New York, NY, USA, 325–329. <https://doi.org/10.1145/3406522.3446046>
- [55] Ben Poole, Ajay Jain, Jonathan T. Barron, and Ben Mildenhall. 2022. DreamFusion: Text-to-3D using 2D Diffusion. *CoRR abs/2209.14988* (2022). <https://doi.org/10.48550/arXiv.2209.14988> arXiv:2209.14988
- [56] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. 2022. Hierarchical Text-Conditional Image Generation with CLIP Latents. *arXiv preprint arXiv:2204.06125* (2022).
- [57] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. 2021. Zero-Shot Text-to-Image Generation. In *Proceedings of the 38th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 139)*, Marina Meila and Tong Zhang (Eds.). PMLR, 8821–8831. <https://proceedings.mlr.press/v139/ramesh21a.html>
- [58] Ali Razavi, Aaron Van den Oord, and Oriol Vinyals. 2019. Generating diverse high-fidelity images with vq-vae-2. *Advances in neural information processing systems* 32 (2019).
- [59] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2022. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 10684–10695.
- [60] Ian Ruthven. 2008. Interactive information retrieval. *Annu. Rev. Inf. Sci. Technol.* 42, 1 (2008), 43–91. <https://doi.org/10.1002/aris.2008.1440420109>
- [61] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S. Sara Mahdavi, Rapha Gontijo Lopes, Tim Salimans, Jonathan Ho, David J Fleet, and Mohammad Norouzi. 2022. Photorealistic Text-to-Image Diffusion Models with Deep Language Understanding. *arXiv preprint arXiv:2205.11487* (2022).
- [62] Tetsuya Sakai. 2007. Alternatives to Bpref. In *SIGIR 2007: Proceedings of the 30th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, Amsterdam, The Netherlands, July 23-27, 2007*, Wessel Kraaij, Arjen P. de Vries, Charles L. A. Clarke, Norbert Fuhr, and Noriko Kando (Eds.). ACM, 71–78.
- [63] Tetsuya Sakai. 2008. Comparing metrics across TREC and NTCIR: The robustness to system bias. In *Proceedings of the 17th ACM Conference on Information and Knowledge Management, CIKM 2008, Napa Valley, California, USA, October 26-30, 2008*, James G. Shanahan, Sihem Amer-Yahia, Ioana Manolescu, Yi Zhang, David A. Evans, Aleksander Kolcz, Key-Sun Choi, and Abdur Chowdhury (Eds.). ACM, 581–590.
- [64] Rob Salkowitz. 2022. Midjourney Founder David Holz On The Impact Of AI On Art, Imagination And The Creative Economy. *Forbes* (Sept. 2022). <https://www.forbes.com/sites/robsalkowitz/2022/09/16/midjourney-founder-david-holz-on-the-impact-of-ai-on-art-imagination-and-the-creative-economy/>
- [65] Greg Schohn and David Cohn. 2000. Less is More: Active Learning with Support Vector Machines. In *Proceedings of the Seventeenth International Conference on Machine Learning (ICML 2000), Stanford University, Stanford, CA, USA, June 29 - July 2, 2000*, Pat Langley (Ed.). Morgan Kaufmann, 839–846.
- [66] Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton Mullis, Aarush Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki. 2021. LAION-400M: Open Dataset of CLIP-Filtered 400 Million Image-Text Pairs. *CoRR abs/2111.02114* (2021). arXiv:2111.02114 <https://arxiv.org/abs/2111.02114>
- [67] Tal Schuster, Roei Schuster, Darsh J. Shah, and Regina Barzilay. 2020. The Limitations of Stylometry for Detecting Machine-Generated Fake News. *Comput. Linguist.* 46, 2 (jun 2020), 499–510. https://doi.org/10.1162/coli_a_00380

- [68] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. 2015. Deep unsupervised learning using nonequilibrium thermodynamics. In *International Conference on Machine Learning*. PMLR, 2256–2265.
- [69] Simon Tong and Edward Y. Chang. 2001. Support vector machine active learning for image retrieval. In *Proceedings of the 9th ACM International Conference on Multimedia 2001, Ottawa, Ontario, Canada, September 30 - October 5, 2001*, Nicolas D. Georganas and Radu Popescu-Zeletin (Eds.). ACM, 107–118. <https://doi.org/10.1145/500141.500159>
- [70] Salahuddin Unar, Xingyuan Wang, Chuan Zhang, and Chunpeng Wang. 2019. Detected text-based image retrieval approach for textual images. *IET Image Process.* 13, 3 (2019), 515–521. <https://doi.org/10.1049/iet-ipr.2018.5277>
- [71] Ellen M. Voorhees. 2001. The Philosophy of Information Retrieval Evaluation. In *Evaluation of Cross-Language Information Retrieval Systems, Second Workshop of the Cross-Language Evaluation Forum, CLEF 2001, Darmstadt, Germany, September 3-4, 2001, Revised Papers (Lecture Notes in Computer Science, Vol. 2406)*, Carol Peters, Martin Braschler, Julio Gonzalo, and Michael Kluck (Eds.). Springer, 355–370.
- [72] Ellen M. Voorhees. 2019. The Evolution of Cranfield. In *Information Retrieval Evaluation in a Changing World - Lessons Learned from 20 Years of CLEF*, Nicola Ferro and Carol Peters (Eds.). The Information Retrieval Series, Vol. 41. Springer, 45–69.
- [73] Ellen M. Voorhees, Ian Soboroff, and Jimmy Lin. 2022. Can Old TREC Collections Reliably Evaluate Modern Neural Retrieval Models? *CoRR abs/2201.11086* (2022). arXiv:2201.11086
- [74] Xintao Wang, Yu Li, Honglun Zhang, and Ying Shan. 2021. Towards Real-World Blind Face Restoration With Generative Facial Prior. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2021*. Computer Vision Foundation / IEEE, 9168–9178. <https://doi.org/10.1109/CVPR46437.2021.00905>
- [75] Ryen W. White, Mikhail Bilenko, and Silviu Cucerzan. 2007. Studying the use of popular destinations to enhance web search interaction. In *SIGIR 2007: Proceedings of the 30th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, Amsterdam, The Netherlands, July 23-27, 2007*, Wessel Kraaij, Arjen P. de Vries, Charles L. A. Clarke, Norbert Fuhr, and Noriko Kando (Eds.). ACM, 159–166. <https://doi.org/10.1145/1277741.1277771>
- [76] Zuobing Xu, Ram Akella, and Yi Zhang. 2007. Incorporating Diversity and Density in Active Learning for Relevance Feedback. In *Advances in Information Retrieval, 29th European Conference on IR Research, ECIR 2007, Rome, Italy, April 2-5, 2007, Proceedings (Lecture Notes in Computer Science, Vol. 4425)*, Giambattista Amati, Claudio Carpineto, and Giovanni Romano (Eds.). Springer, 246–257. https://doi.org/10.1007/978-3-540-71496-5_24
- [77] Christoph Zauner. 2010. *Implementation and benchmarking of perceptual image hash functions*. Master’s thesis. Upper Austria University of Applied Sciences, Hagenberg Campus.
- [78] Rowan Zellers, Ari Holtzman, Hannah Rashkin, Yonatan Bisk, Ali Farhadi, Franziska Roesner, and Yejin Choi. 2019. Defending Against Neural Fake News. In *Advances in Neural Information Processing Systems*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett (Eds.), Vol. 32. Curran Associates, Inc. <https://proceedings.neurips.cc/paper/2019/file/3e9f0fc9b2f89e043bc6233994dfcf76-Paper.pdf>
- [79] Cha Zhang and Tsuhan Chen. 2002. An active learning framework for content-based information retrieval. *IEEE Transactions on Multimedia* 4, 2 (2002), 260–268. <https://doi.org/10.1109/TMM.2002.1017738>
- [80] Yinglong Zhang and Robert Capra. 2019. Understanding How People Use Search to Support Their Everyday Creative Tasks. In *Proceedings of the 2019 Conference on Human Information Interaction and Retrieval* (Glasgow, Scotland UK). ACM, New York, NY, USA, 153–162. <https://doi.org/10.1145/3295750.3298936>
- [81] Yinglong Zhang, Rob Capra, and Yuan Li. 2020. An In-Situ Study of Information Needs in Design-Related Creative Projects. In *Proceedings of the 2020 Conference on Human Information Interaction and Retrieval* (Vancouver BC, Canada). ACM, New York, NY, USA, 113–123. <https://doi.org/10.1145/3343413.3377973>