

Evaluating Digital Library Search Systems by using Formal Process Modelling

Christin Katharina Kreutz*, Martin Blum†, Philipp Schaer*, Ralf Schenkel‡, Benjamin Weyers‡

*TH Köln - University of Applied Sciences, Cologne, Germany

†Schloss Dagstuhl LZI, dblp, Trier, Germany

‡Trier University, Trier, Germany

*{christin.kreutz, philipp.schaer}@th-koeln.de, †martin.blum@dagstuhl.de, ‡{schenkel, weyers}@uni-trier.de

Abstract—Evaluations of digital library information systems are typically centred on users correctly, efficiently, and quickly performing predefined tasks. Additionally, users generally enjoy working with the evaluated system, and completed questionnaires show an interface’s excellent user experience. However, such evaluations do not explicitly consider comparing or connecting user-specific information-seeking behaviour with digital library system capabilities and thus overlook actual user needs or further system requirements.

We aim to close this gap by introducing the usage of formalisations of users’ task conduction strategies to compare their information needs with the capabilities of such information systems. We observe users’ strategies in scope of expert finding and paper search. We propose and investigate using the business process model notation to formalise task conduction strategies and the SchenQL digital library interface as an example system. We conduct interviews in a qualitative evaluation with 13 participants from various backgrounds from which we derive models.

We discovered that the formalisations are suitable and helpful to mirror the strategies back to users and to compare users’ ideal task conductions with capabilities of information systems. We conclude using formal models for qualitative digital library studies being a suitable mean to identify current limitations and depict users’ task conduction strategies. Our published dataset containing the evaluation data can be reused to investigate other digital library systems’ fit for depicting users’ ideal task solutions.

Index Terms—digital libraries, information-seeking behaviour, user study, qualitative evaluation, human focus

I. INTRODUCTION

Nowadays, there are numerous options for obtaining bibliographic information. Classic bibliographic digital libraries (DLs) such as the ACM DL, Bibsonomy, dblp, Google Scholar, Semantic Scholar, SpringerLink, or ResearchGate support search and exploration functionality. More complex information needs that do not solely require using keyword-based search can be assessed by more sophisticated information systems. One option is to write Cypher queries in the specialised tool GrapAL [1] or to formulate RDF-style triples in the Narrative Service of PubPharm [2]. Another could be using GUI-assisted query construction such as SchenQL [3].

Evaluations of such DL systems are crucial for their development as they help identify current shortcomings and room for improvement [4]. Historically, library and information science focused more on systems than on users’ perspectives [5]. Nowadays, user studies are conducted in many forms

with a plethora of different measures [4] in the domain of bibliographic metadata: Many approaches assess quantifiable measures such as the correctness of user-constructed queries to fulfil specific tasks [6, 7, 8], time spent to satisfy an information need [6, 7, 8], the constructed query size [8, 9], the number of clicks to find the solution for a task [7] or questionnaires with a focus on deriving subjective user feedback [6, 7, 10] which may contain quantitative scales, such as Likert scales. Some DLs are evaluated more qualitatively by using think-aloud protocols [7, 11, 12], query log analysis [9], open-ended questions for users [13, 14, 15, 16] or interviews with domain experts [1, 6, 17].

However, Kuhlthau [5] argued for the need to connect users’ information-seeking behaviour and the systems providing the information, to think beyond precision-based evaluations and include users’ perspectives in the search process. This corresponds to Kelly’s [18] definition of *Information-Seeking Behavior with IR Systems* as human-focused studies investigating users’ usual information-seeking behaviour while interacting with an information system. Her mentioned prototypical studies of this category [19, 20, 21, 22] all produce quantifiable data that can be statistically analysed.

Overall, current DL evaluation designs and studies appear to be primarily concerned with quantitative measures, with the occasional report of anecdotal insight gained from free-text questions or (semi-structured) interviews. Considering this quantitative data may give hints to which elements of a DL may generate specific issues reducing search performance, it lacks of giving clear directions of further improvements. In this case, qualitative data with a broader scope could give better insight into the actual user needs as well as further indications of potential future developments. This is also the primary function of qualitative methods in designing interactive systems [23]. A crucial factor for the design of interactive systems, including DLs, is the actual task conducted by the user with the system [24]. With this work at hand, we aim to consider the actual task conduction model of such search tasks conducted in DLs as the central element for the successful design and evaluation of DLs in this context [4]. Still, despite studies observing usage models [25], a specific qualitative evaluation of discrepancies between users’ ideal task execution strategies and a systems’ design has been disregarded so far.

The lack of these studies might be attributed to difficulties

in formalising information exploration processes with standard process modelling tools [26]. Exploratory search is frequently preceded by an anomalous state of knowing where a searcher requires information to proceed in their task [27]. So the formulation of information needs is an iterative process in which new information influences a user’s perception of the information space [27]. Quantifying these changes in a user’s task conduction model could assess changes in the information formulation [27] as a fine-grained micro-observation. Contrasting this viewpoint, our macro-objective is capturing the difference in the overall representation of the task execution process of users and the functionality offered by an information system. We do not consider single conceptual changes of the users’ task conduction models during information exploration. Instead, we investigate the users’ task conduction models of the complete task in a cumulative, outcome-oriented form.

Therefore, in this work, we formalise the task conduction models derived from users’ typical task completion strategies to compare and connect them with the capabilities of an information system. We propose investigating how users’ task conduction models are translated to process documentations and adapted to fit the limited options of a single system with the following general research question: *How can we compare users’ conceptions of search tasks in digital library with capabilities of such a system?* As a goal, we seek to help identify dissonances in an interface’s functionality that could be modified to suit the users’ pre-existing conceptions of typical tasks in such systems better

As an exemplary application of the proposed method, we use the SchenQL ecosystem [3] in our study. At the time of writing this paper, this is the most recent digital library interface suitable to satisfy many complex information needs, surpassing the scope of other current DLs [3]. SchenQL’s usage is easy to learn, and the system is appropriate for users of DLs with different expertise levels [6]. We use Law et al.’s [28] BPMN variant to represent users’ task conduction models, as this method does not incorporate a modellers perspective but instead solely focuses on capturing a user’s perspective collected through interview data. In our *Information-Seeking Behavior with IR Systems* type study (i.e., [29]), we analyse and compare BPMNs of users’ ideal task conduction [28], models describing users’ actual task conduction using a specific system (we use SchenQL [3] as an example for this work) and models translating the initial system-independent strategy to the system (see Figure 3).

Our contribution can be summed up as follows:

- Application of a formal representation method of user processes [28] into a digital library evaluation setup with a specific focus on identifying current shortcomings.
- Description and analysis of users’ ideal task conduction models for everyday tasks in DLs as well as discrepancies and adaptation when confronted with a specific DL.
- Exemplary analysis of our qualitative evaluation method on the SchenQL query language and user interface.
- Publication of reusable interviews and formal models of users’ processes of solving typical tasks in DLs [30].

II. RELATED WORK

We discuss areas adjacent to this work: We present *bibliographic digital libraries* before compiling strategies for the *evaluation of bibliographic digital libraries* and presenting different types of *modelling users’ task conduction*.

A. Bibliographic Digital Libraries

The ACM DL, Bibsonomy [31], dblp [32], Google Scholar, Semantic Scholar [33], Springer Link, ResearchGate, Clarivate/Web of Science [34], Elsevier Scopus [34], and Dimensions [35] are DLs operating on bibliographic metadata which offer keyword-based search to retrieve and explore the underlying data. Some systems in the domain of bibliographic metadata provide more functionality or information than simple keyword-based search: OpenAlex [36] provides a web API and database snapshot, which allows the formulation of complex queries. Semantic Scholar [33] extracts paper content information, such as entities and mentions, which are used to construct a literature graph. Clarivate/Web of Science [34] offers a sophisticated faceted search interface to construct queries. GrapAL [1] is a Cypher-based query formulation tool and GUI offering complex graph operations. Daffodil [11] is a faceted digital library interface that supports users in expressing their complex information needs by suggesting query components and marking potential errors. SchenQL [3] (see Sec. III-B) is a domain-specific query language and GUI supporting information search and exploration by providing aggregations and expert functions.

B. Evaluation of Bibliographic DLs

Opposing our qualitative evaluation of a DL, multiple *quantitative* study designs have been conducted, mainly focusing on users solving predefined tasks [6, 7, 8, 9, 11]. Quantitative measures such as answering time [6, 8], correctness [6, 8], query size [8, 9], used query components [9], perceived query difficulty [6] and query execution time [6] were reported.

There are also some *less structured evaluations* focusing on finding challenges of current systems, and users’ information needs through expert interviews [1, 6, 17], usage diaries [25], questionnaires [10, 12, 13, 15, 16, 25], user assessment of prototypical interfaces [11, 17], thinking aloud interviews while conducting tasks [11] and screen captured task conduction which was later annotated by users [12]. While these methods targeted one of the goals of our study, they did not systematically connect users’ perspectives with the respective systems. From thinking-aloud tests [11], formal cognitive models could have been derived. However, a systematic analysis of the task conduction and a comparison using the participant’s usual DLs was disregarded, compared to our approach. Furthermore, when evaluating usage diaries [25], no tasks were predefined, contrasting our approach; another difference is disregarding the user’s typical task conduction strategy without the system at hand. Another evaluation [12] only observed the typical task execution behaviour of participants and constructed formal models from the transcribed data but did not relate it back at a single system.

C. Modelling Users' Task Conduction

Task modelling has been developed in the context of engineering of interactive systems. This family of model-based approaches aims at informing and driving the design and development of interactive systems [24]. Either user or system prospective could be taken where the user perspective is relevant for our work. The driving question here is how a user performs a search task in a given DL such that task models enable to describe these processes and reuse them in terms of creation of DL interactive interfaces.

Various modelling approaches and languages have been proposed in the area of engineering of interactive systems. ConcurTaskTrees (CTT) [37] is an approach to model all actions users partake to achieve specific goals. The notation can be used to describe how tasks could be solved in an existing or envisioned system or how users think a task should be performed. CTT focuses on activities, has a hierarchical structure, a graphical syntax and provides numerous temporal operators. HAMSTERS [38] is a notation based on and compatible with CTT used to decompose complex real-life models into smaller task-based, interconnected ones. Task models can communicate with each other. The notation distinguishes between task types such as system, users (e.g., cognitive or motor tasks) and interactive tasks (e.g., input or output tasks). Temporal relationships between tasks are important.

Dias et al. [26] present the Exploratory Search KiP model to visualise information-seeking in the web. They model four activities: Search term selection, query formulation, result check and information extraction. From click tables and explanations for clicks derived from thinking-aloud protocols, these activities are defined [39]. The modeller is explicitly involved in the design process by incorporating their interpretation of users' actions and statements to identify mismatches.

Still, literature around the creation of task models mostly neglects the actual creation process as has been identified by Bowen et al. [24]. With the Business Process Modeling Notation (BPMN) [40], processes can be formalised in detail. Processes can be depicted through sequences of activities, conditions, paths and logical operators. The variant introduced by Law et al. [28] consists of a subset of the BPMN modelling options for depicting structured user interaction models from thinking-aloud interviews. They proposed a methodology to gather BPMN models in a structured way from an informal data base. These have been gathered from so called think aloud interviews, a qualitative method combining think aloud protocol with open interview techniques (see Section III-A).

III. METHODOLOGICAL BASICS

This section describes the combined methods in our evaluation, briefly introduced in the previous section.

A. Modelling Users' Task Conduction

Law et al. [28] describe a sequential and imperative method to construct structured user interaction models from unstructured, process-oriented thinking-aloud interviews. They create models in BPMN by segmenting transcribed interviews into

tasks and later classifying these segments into six categories. The process focuses on events, tasks and conditions as provided by BPMN. Artefacts such as tools or data are not specifically modelled.

Segments are created by cutting the transcripts by verbs, as they signal tasks or events. For adjective or relative clauses, the segments should not be separated. Time-related phrases should be separated, even without the presence of verbs. In general, a finer segmentation is preferred.

In the following classification step, segments are identified to be either of type `setting`, `annotation`, `task`, `event`, `condition` or `other`. In our use case, the most prominent classes are tasks (activities performed by users; yellow parts in Fig. 6), events (events occurring during task conduction, incidents in the environment; red parts in Fig. 6), conditions (descriptions of alternatives; blue part in Fig. 4) and annotations (mentioned DLs, systems, observations from screen captured data; green part in Fig. 4). These classes correspond to specific BPMN elements. Figure 1 shows an example transformation of an interview to a (partial) BPMN.

This method does not focus on the details of activities but rather on switching tasks. Its structure aims at preventing including the modeller's interpretation of the interviews and focuses on the interview content. This has been demonstrated in the original article by an empirical evaluation of the method.

B. The SchenQL Query Language and GUI

SchenQL [3] is a domain-specific query language and graphical user interface (GUI) supporting information search and exploration. Its main goal is to support various types of users of digital libraries in their information needs. It offers a broad selection of domain-specific functions such as co-author search, citation aggregation or providing bibliographic metrics. The system offers query formulation following a predefined sophisticated grammar as well as information exploration via the interface. Its GUI offers suggestions and auto-completion of query components to help users with query formulation (shown in Figure 2). SchenQL was designed with the dblp computer science bibliography as the central use case. Kreutz et al. [3] give an overview of SchenQL's grammar and GUI.

We selected SchenQL as an example DL in our studies for its recency, meaning the system is still being in use or updated [3], its suitability for users of different experience with DLs [6], its ease of use [6] and mainly for its broad coverage of bibliographic information needs, surpassing the scope of most current DLs [3].

C. General Idea

Our general idea is to use users' usual process for task conduction and their solution strategy when using a specific digital library system to evaluate DL systems. For this, we compose three models for each task and user: *First* we construct a BPMN from a study participant's interview on their general process for solving a task using their preferred digital libraries following the method by Law et al. [28]. Then, the same user verifies their ideal task conduction model, which

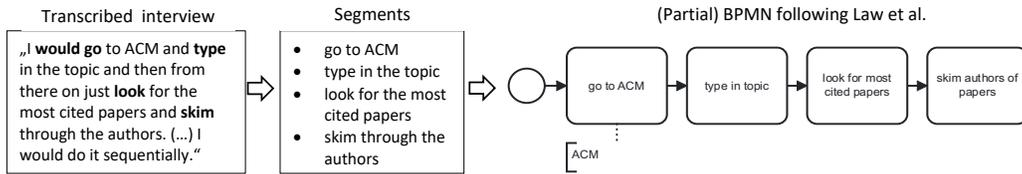


Fig. 1. Construction of the (partial and unverified) BPMN from transcribed interview data for participant *blue_dog*.



Fig. 2. SchenQL search interface with colour-coded (dependent on the component type) query component suggestions.



Fig. 3. Construction of the three task conduction models: The vIMM comes from a participant's general usage of digital libraries or other sources, the vPGM comes from an expert's translation of the vIMM to the DL which is being evaluated and the PCM comes from the participant using that system.

results in their verified **Interview Mental Model (vIMM)**. The *second* model we construct (the **verified Process Gold Model, vPGM**) is an expert-generated translation of a user's workflow from their chosen DLs and tools to the example digital library system. This BPMN represents the application of a user's strategy as closely as possible to the DL being evaluated. The *third* model (the **Process Conduction Model, PCM**) is composed by a user's actual execution of the task with the DL, which is being evaluated. This BPMN represents users' actual task solution process using a specific system.

Figure 3 shows a simplified construction of the three models. In this work, we choose SchenQL [3] as an exemplary DL to demonstrate the evaluation process and analyse the results.

IV. EXPERIMENTAL SETUP

This section presents our research questions, tasks, the study participants, our evaluation setup, and technical details.

We compare the model of users' general task conduction produced with Law et al.'s [28] method with a model composed from users' task conduction in a specific system (SchenQL [3]) and a model of the translation of their initial system-independent strategy to the system (see Figure 3).

A. Research Questions

We observe different aspects related to our general question *How can we compare users' conceptions of search tasks in a digital library with capabilities of such a system?* We investigate the following five more fine-grained research questions

by using Law et al.'s [28] method to represent users' task conduction models and SchenQL [3] as the example DL system for this work to observe the suitability and use of formalised task conduction models:

- RQ₁ What are users' preferences, which components of digital libraries are usually used for the predefined tasks?
- RQ₂ How do users utilise the example system, which components are used for the specific predefined tasks?
- RQ₃ What are the limitations of the example DL system? Which components or functions were ignored or missed?
- RQ₄ Is the example system usable for advanced DL tasks?
- RQ₅ What are the discrepancies between the ideal task conduction models of users and their actual task conduction, how are models adapted to solve the predefined tasks?

B. Tasks

We focus on two deliberately vague exploratory task descriptions, giving study participants some wiggle room to fit their usual information-seeking behaviour better. For instance, a user might define *expertise* as having published many papers in specific topic-related journals, while another might consider those with a high topic-independent *h* index as relevant.

- Task T_{ex} : Find two experts on a topic of your liking.
- Task T_{pa} : Find relevant papers from a topic of your liking which appeared after 2017.¹

¹Soufan et al. [41] provide an overview on the exploratory search task.

C. Participants

Our thirteen study participants are computer or information scientists with differing expertise in using DLs for research tasks, by which they were invited to take part: Two masters students, six PhD students (first year to last year students), an industry researcher, a dblp staff member, a postdoc and two professors. They participated voluntarily and did not receive any incentives. For anonymisation, participants chose code names with which we refer to them.

D. Setup and Experiment Conduction

Our experimental setup is composed of four consecutive parts, with seven steps. Users were part of two user sessions; other parts did not include the study participants. They were evaluated on a one-by-one basis with one (the same) investigator present. In the following, the seven steps are described.

Step i): Pre-Interview Questionnaire. Step *i*) and *ii*) were conducted sequentially and took about 30 minutes. They form the first user session. Participants filled out a questionnaire to consent to voluntarily participate in the study, have their voices and screens recorded, and have transcripts as well as formalised task conduction models published. All participants were explicitly made aware that they could stop and drop out of the experiment at any time without consequences. The process (a mail to the investigator) for later deleting user-specific data from the resulting dataset was also explained. The conduction of an initial questionnaire was accompanied by an interviewer and audio recorded.

Step ii): Interview. In recorded semi-structured 1:1 interviews with the participants, they described how they usually conducted tasks T_{ex} and T_{pa} with a focus on the data sources. For T_{ex} , the participants' topics, their familiarity with the topics, their definition of an expert and their process of solving T_{ex} were asked. Afterwards, the participants chose a topic for the second task, indicated their expertise in this topic, and gave their notion of relevancy before they described their process to solve T_{pa} . The interviewer did not intervene with the participants' understanding of the tasks but posed clarifying questions regarding temporal or process-related ambiguities.

Step iii): Modelling I. From transcripts of the interviews for participants, one modeller constructed mental models, so-called *Interview Mental Models (IMMs)*, for all tasks and users. The formation of IMMs with BPMN followed the description by Law et al. [28] with mentioned information sources such as digital libraries included in annotations.

Step iv): Verification. Step *iv*) to *vi*) (the second user session) were done in one session and took about one hour. This step verified the modeller-generated IMMs by the participant. Finally, the investigator went through the IMMs for the two tasks and asked the participant to intervene and note changes if the formalisation did not depict that user's usual process.

Step v): Tasks. Participants watched the five-minute video demonstration² of SchenQL by Kreutz et al. [3] and had access to a language documentation. The screen (e.g., users' mouse

movement or entered queries) and think-aloud protocol were recorded. Users conducted T_{ex} and then T_{pa} and stopped with a task when feeling they solved it or after 25 minutes for T_{ex} and 15 minutes for T_{pa} .

Step vi): Post-Task Questionnaire. In a questionnaire, participants indicated the components that they enjoyed using in the SchenQL interface, components missing or not found and encountered problems. They could also mention anything they did not get the chance to or forgot to mention earlier.

Step vii): Modelling II. First, the modeller revised the IMMs according to the annotations made in step *iv*) to represent the user's information-seeking strategy. The resulting model is the so-called *verified IMM (vIMM)*.

Based on vIMMs, a SchenQL expert translated the workflow as closely as possible to the SchenQL system in a *Process Gold Model (PGM)* in BPMN. SchenQL queries were included as annotations. A second independent SchenQL expert who was previously uninvolved in the evaluation process verified the models. When the two SchenQL experts disagreed on the workflow, they discussed the issue until reaching a consensus. The *verified PGMs (vPGMs)* depict the revised translation of the vIMMs to the example DL SchenQL.

For all users and tasks, a user's actual *Process Conduction Model (PCM)* was constructed from their think-aloud protocol with BPMN following Law et al. [28] when using the exemplary system. Screen-captured data and action non-inducing verbs such as *assume*, *feel*³ and *wonder* were used to annotate the PCM, e.g., with constructed SchenQL queries. Non-descriptive sentences such as "*oh, interesting*" while clicking a result or "*okay, h index of 40*" are used for further annotation. When a participant was unable to finish a task in the reserved time frame, the PGM includes the information of the task conduction stopping.

Our experiment produces three models for each user and task: A user's verified general task conduction (vIMM), a verified translation of the task to the single digital library to evaluate (vPGM), and a model of the user's actual task conduction using that DL (PCM).

E. Technical Details

As data source, we follow Kreutz et al. [3] and use the dblp XML dump from 1st Oct. '21⁴ and Semantic Scholar data [33] from Oct. '21 for citations and AMiner Open Academic Graph 2.1 [42, 43, 44] for identifying automatically generated keywords of publications; abstracts are taken from both collections. We used Descript⁵ for transcription and manually corrected the transcripts. We follow the transcription rules of Dresing and Pehl [45]. De-anonymising parts of the transcripts (e.g., mention of names) and the questionnaires were marked and replaced. BPMNs were created with BPMN.io⁶.

³The classification of segments which describe feelings as *annotation* deviates from Law et al.'s [28] to improve consistency of BPMNs. Such sentences normally explain user behaviour. In our case, most explanations stem from screen capture observations and are modeled as *annotations*.

⁴<https://dblp.org/xml/release/dblp-2021-10-01.xml.gz>

⁵<https://www.descript.com>

⁶<https://bpmn.io>

²<https://youtu.be/pkaKe7vo9ys>

V. OBSERVATIONS FROM COLLECTED DATA

The following observations stem from the data we collected during our study via interviews and the constructed formal models for the two tasks of our thirteen study participants.

First, we present and discuss the three models of a specific participant as an example before the general and more fine-grained observations are laid out.

A. Exemplary Data Series in Context of all Data

To highlight the differences in constructed models and the general reoccurring observations or problems, this section discusses the BPMNs of participant *green_deer*⁷ for T_{pa} concerning all observed data. Figures 4, 5 and 6 depict the vIMM, vPGM and PCM.

In general, several (4 for T_{ex} , 7 for T_{pa}) vIMMs contain multiple arms to begin with, which depend on a condition or are usually conducted in parallel (see, e.g., *actively searching*, Fig. 4). For some vIMMs we encountered the problem of some parts of models being very unspecific as seen with the segment *check authors' other works' importance* (see Fig. 4). Many participants did not mention with the support of which specific system they conducted some segments (see, e.g., *check #citations*, Fig. 4). Several segments could not be translated to SchenQL at all, e.g., if they included outside help (see, e.g., *check mouth propaganda*, Fig. 4), or needed to be modified in the translation process (see, e.g., *go to arXiv*, Fig. 4 translated to parallel tasks in Fig. 5). With PCMs we often faced the problem of participants not describing what they were doing or intending to do. They also rarely followed their usual strategy described in their vIMM. PCMs thus heavily rely on observations from screen capture such as clicking information or entering queries. While users conducted tasks in SchenQL, the investigator did not interfere with their understanding of having solved the task, so the PCMs depict the user's satisfaction rather than the correctness of their result.

More specific to the conduction models for *green_deer*, the participant seemed to only focus on papers about their chosen topic, they did not search for (known) existing implementations on the topic (compare Fig. 4 and 6). The source or novelty of resulting publications also seemed to be irrelevant when using SchenQL as not only preprints were searched for (compare Fig. 4 and 6). This lack of restricting the result set might stem from the user's limited knowledge of SchenQL's capabilities. Comparing the vIMM with the PCM leads to the observation that the participant only actively tries⁸ to incorporate one (*check #citations*, see Fig. 4) of the four described decision criteria. From our collected data we cannot determine if they, e.g., also included the check of authors' other works' importance by recognising resulting paper titles and remembering their authors and their full body of work.

Green_deer seemed only to do part of their vIMM to solve the task. This specific part can be recognised in the PGM, where it has been fleshed out in more detail.

⁷The participant explicitly consented to their BPMNs being used as example models in this paper.

⁸From T_{ex} we assume they confuse *references* with *citations*.

B. Verified Interview Mental Models (vIMMs)

1) *vIMMs for T_{ex}* : Popular components in solving T_{ex} are keyword search (12) for papers or venues (described by Bates [46] as identifying central venues), considering authors of popular or somehow good papers as experts (7), looking at the number of citations (7), affiliations (6) and taking a deeper look into the paper, e.g., by checking the references in the related work or introduction (6). Five participants incorporated some undefined notion of relevancy or expertise. Three participants searched for surveys, looked at abstracts of papers, considered author positions on papers, venue ranking, asked others for their opinion or followed references of papers (while only one participant followed a paper's incoming citations). Related terms and query formulation or the most cited papers were important for two. All participants, except one, usually use digital libraries to solve the expert search task.

For this task, only two participants described using a single system, nine users explicitly described a switch of their used tool throughout their process. Only two participants described the process with multiple (disjunctive) options to start solving the task. As for the system or tool to start with, multiple mentions were possible: Four described using Google, the same amount of people started the process by using Google Scholar. Throughout their process, over half (7) of the participants used Google Scholar, and the same amount incorporated a Google search. These observations partially mirror users' previous indications of DLs or systems they normally use.

2) *vIMMs for T_{pa}* : Most participants (12) would use keyword search in T_{pa} . Less than half of the participants (6) incorporate following references (described by Bates [46] as backwards chasing), using related terms or refining queries and asking other people for their opinion or expertise while searching for relevant papers. In five cases, abstracts of papers are read or participants take a deeper dive into papers, e.g., check references in the related work section or look at figures of papers. The same amount of people use a somewhat unclear definition of relevance such as "*author's other works' importance*" or "*good authors*". The number of citations is relevant for four, venue rankings and following citations are part of three vIMMs. Constructing a citation graph and searching for surveys were parts of two models each.

For this task, no participants described using only a single system, ten users explicitly described a switch of their used tool throughout their process. Five participants described the process with multiple (disjunctive) options to start solving the task. As for the system or tool to start with, we did not encounter clear preferences: Four described using Google, three people started the process by using Google Scholar and dblp. Throughout their process, only five participants described using Google Scholar, four mentioned using Google. Used systems seemed to vary more compared to T_{ex} .

C. Translation of Verified Interview Mental Models (vIMMs) into Verified Process Gold Models (vPGMs)

We encountered several problems with translating vIMMs to vPGMs with our example DL. We were not able to formulate

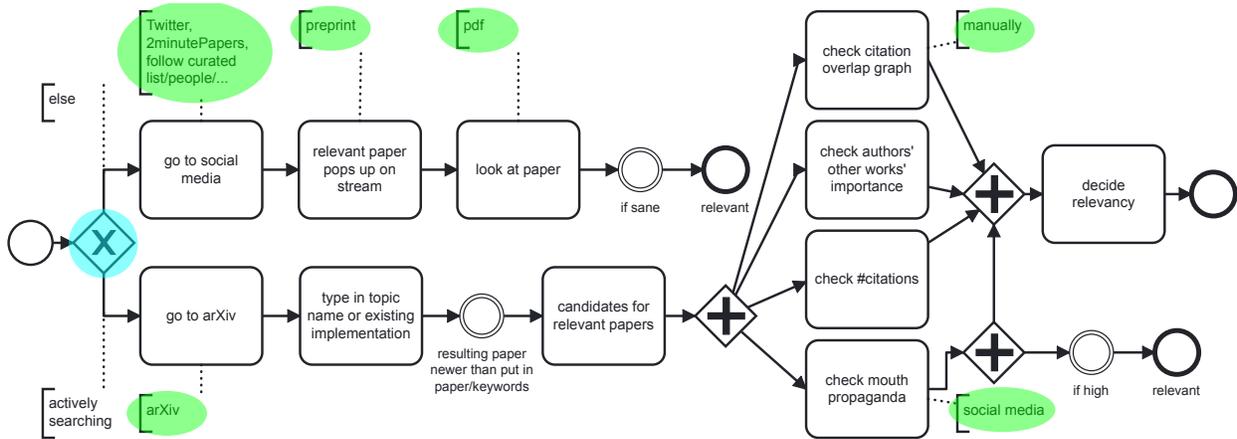


Fig. 4. vIMM of T_{pa} for participant *green_deer*.

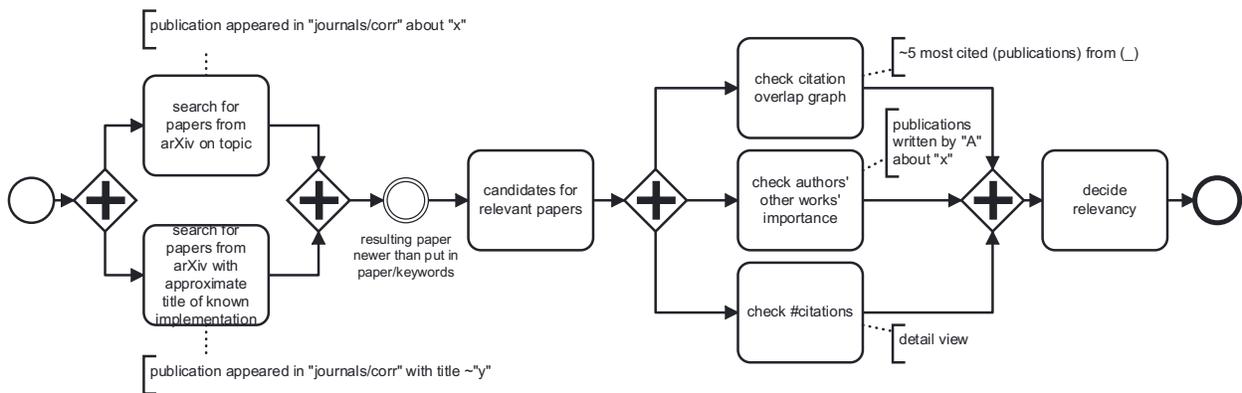


Fig. 5. vPGM of T_{pa} for participant *green_deer*.

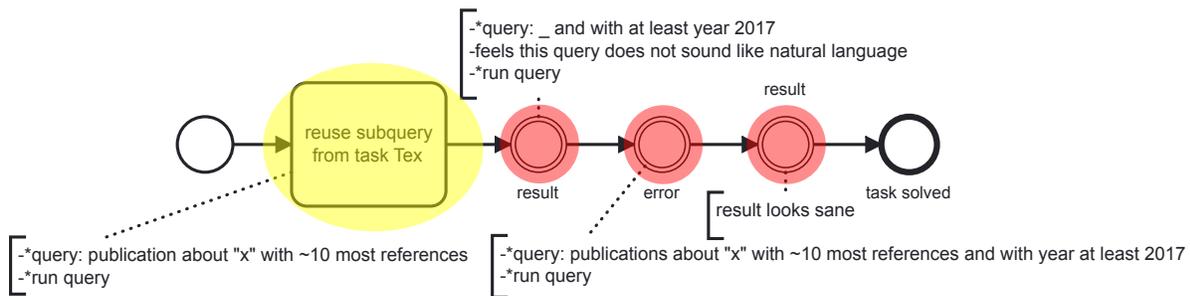


Fig. 6. PCM of T_{pa} for participant *green_deer* with annotations preceded by * being observations from the screen capture.

several segments of users' usual workflow throughout both tasks: Eight cases of participants getting help from humans in their vIMMs in T_{ex} , five times someone used a general keyword-based search in Google to get an overview of their topic or trust Google's ranking, two users wanted to incorporate author positions in papers into their workflow, four instances occurred where persons intended to check papers related to a specific one. Twice the restriction of papers from a specific publisher was not modelled and one participant wanted to contact persons directly. In two cases, we encour-

tered used data sources or information not contained in our dataset. In three cases, we logically reordered segments in vPGMs compared to vIMMs. All these cases were related to downloads of pdfs and checking the information in these pdfs. We could eliminate the restriction of asking others for information once, as the underlying dataset contained the data. In another case, we were able to support performing a segment for which information was missing if a user was less familiar with a topic. One vIMM was not translated to a vPGM as it contained no components that could be performed in SchenQL.

D. Process Conduction Models (PCMs)

1) *PCMs for T_{ex}* : All participants queried for some form of publications about their specific topic with a publications about type query. Eight persons even used this query as their initial one. Many participants (9) looked at the example queries in the GUI; five even used the documentation for an in-depth look. Even though T_{ex} is about persons, only five users checked any person’s profile, while most (9) opened detail views for publications. Slightly more participants (6) constructed a query that asked for persons (persons authored). Four users looked at co-authors of people. Abstracts of papers were read in five cases, BowTie visualisations [47] were actively used slightly less (4). As for query components, four participants used the aggregation function with, the same amount queried for most cited and three used \sim RANK.

2) *PCMs for T_{pa}* : For this task, participants did not need to consult the example queries (5) and documentation (4) as much as for the first one. Almost everyone’s (12) first run query is some form of publications about a topic. Ten persons incorporated a with year component in one of their queries. Slightly fewer participants (9) took the time to check the publication detail view, only six read an abstract. Used query components of five participants were the aggregation with with and \sim RANK. Three participants used most cited. Two persons each followed citations and references, checked the full-text of papers, specifically looked at authors of papers or re-watched part of the demo video.

3) *Limitations and Problems*: For T_{ex} participants encountered more problems than for T_{pa} . The following tuples depict the number of participants with the specific problem in the task conduction, e.g., (2, 3) describes two participants having trouble in T_{ex} while three users encountered this problem in conducting T_{pa} . Participants encountered several errors with the SchenQL system: They faced errors without indicating their error, which lead to query reformulations (9, 7). Some were met with a chip replace error⁹ in their query where they clicked a suggestion and a query component was replaced by the chip (4, 1), one person was hindered by an error in an example query. Some participants struggled due to the limited scope of the underlying data source (2, 1).

Non-system-caused problems that participants encountered were the following: Uncertainty about the ordering of results (7, 0), insufficient displayed information in results for queries (4, 1), confusion about the difference between full-text search and not (4, 1), complicated syntax (2, 0), confusion about the rank operation (2, 3) and confusion about the difference of terms, keywords and strings (2, 1). Two participants started by using the system like Google. The following problems occurred for single users: Uncertainty about comparing numbers of citations, struggling to filter by the number of citations, confusing citations with references, uncertainty about the difference between most cited and cited by, spelling

⁹Bug, where a user clicking a suggested query chip to continue their query leads to the previous word of the query being wrongfully replaced by it.

errors in their topic, trouble distinguishing colours, confusion on when to use and and with, uncertainty which year at least refers to as starting year and problems with the combination of query parts.

E. Post-Task Questionnaire

1) *Disliked and Desired Components*: Participants mentioned the following components which they disliked: The complicated syntax (4), non-descriptive errors (3), the auto-completion hung up (2), long loading times (2), the unclear sorting (2), the unintuitive system (1), the BowTie (1), complicated query adaption (1), need for brackets (1), the colouring (1) and the map of institutions (1). Participants wished for the following components: Mouseover descriptions or previews (3), more information in the GUI (3), more example queries (2), a documentation that is better searchable (2), saving option for partial queries (1), suggestions from examples (1), dynamic linting (1) and syntax correction options (1).

2) *Liked Components*: Study participants mentioned they liked the auto-completion feature (2) of SchenQL, its query component suggestion (3) and the BowTie visualisation (3) for references and citations of, e.g., papers and persons. Several (6) mentioned the example queries to be very beneficial.

VI. RESULTS

This section analyses the collected questionnaire data and the constructed task conduction models (see Sec. V) for our thirteen participants concerning the aforementioned (see Sec. IV-A) research questions.

A. RQ_1 : User’s Preferences in Digital Libraries

RQ_1 : What are users’ preferences, which components of digital libraries are usually used for the predefined tasks? For this RQ, we focus on analysing the vIMMs concerning the conducted tasks, used systems and reoccurring patterns to find out which parts are usually used for our two tasks.

The exploration strategies for participants seem to be very diverse and consist of complex parallel steps. We see a tendency to start with a search engine when participants describe their usual task conduction for T_{ex} , for T_{pa} we encounter less clear tendencies. Generally, the second task’s strategies seem to rely on a more diverse set of tools.

Almost all participants described using more than one specific system and specifically incorporating a change of used tools in their process, as used systems do not seem to support all their information-seeking strategies. This might lead to switching the working sphere [48], which could negatively affect a user’s cognitive load [49] as well as the degree of user interaction and the systems’ usability [50].

The constructed models for T_{pa} seem to be less complex than the ones for T_{ex} , which could be interpreted as T_{pa} being easier than T_{ex} . The observed need for explicitly changing tools in the second task seems to contradict this assumption.

There are numerous vaguely defined components in the vIMMs which can be interpreted diversely (e.g., *check authors’ other works’ importance* as seen in Figure 4). This phenomenon is a characteristic of users’ pre-focus stage [51].

Some participants told us about another participant having been their role model/mentor in composing their exploration strategy (information on the pairs could threaten anonymity so we do not report it), but we did not see noteworthy overlap in the pairs' strategies. This leads us to assume that users of DLs adapt their usage strategy depending on their preferences.

B. RQ₂: Used Components

RQ₂: *How do users utilise the example system, which components are used for the specific predefined tasks?* In this research question, we analyse PCMs from study participants' actual task conductions with SchenQL.

Participants of our study heavily relied on the provided example queries to start using the system. In general, participants' tendency to explore the search space by searching for papers on a topic and deciding on experts based on their authored works became apparent. Most users stuck to comparatively simple search queries, and only few incorporated more advanced concepts such as aggregations, venue ranks or citation counts in queries. If such measures were considered, they were checked in detail views rather than queries and only in the expert search task. The feature to investigate a query was only used once.

C. RQ₃: System Limitations

RQ₃: *What are the limitations and missing components of the example system?* To answer this, participants' vPGMs, PCMs and post-task questionnaires were observed.

SchenQL cannot model interaction with experts, which is a characteristic of users' pre-focus stage in information search [51], incorporation of other data sources or tools and the type of explorative search found with Google. The errors SchenQL throws are non-descriptive, participants were uncertain about the ordering of results, they required more data to be displayed and they were confused with the differences in full-text search modes. Study participants disliked the example system's complicated syntax. All problems are described in Sections V-D3 and V-E1.

D. RQ₄: Usability for Advanced Tasks

RQ₄: *Is the example system usable for advanced tasks of users of DLs?* The tasks chosen in our evaluation are ones of users of digital libraries, which can be solved by using SchenQL in theory. This RQ examines if users can satisfy their information needs even if some wished-for components they usually utilise might be missing. All except one vIMM could be translated to vPGMs; this model for T_{ex} did not contain any component which SchenQL supports. Therefore we conclude that although several modifications had to be done (see Sec. V-C), SchenQL is generally usable for conducting users' processes. Seven participants felt they solved T_{ex}, eight participants indicated they solved T_{pa}.

E. RQ₅: Discrepancies and Adaption of Task Models

RQ₅: *What are the discrepancies between the ideal task conduction models of users and their actual task conduction,*

how are models adapted to solve the predefined tasks? Here, one DL expert and one SchenQL expert uninformed in the model construction compare vIMMs/vPGMs and PCMs.

The experts agreed that only one participant followed their vIMM for T_{ex} and two for T_{pa}. They found little overlap in vIMMs and PCMs. Some participants seem to have forgotten their vIMM or do not know how to translate it to SchenQL while trying to resolve syntax errors, while others without many syntax problems (partially) stuck to their vIMMs.

In the beginning, participants seemed to follow their models; then they were simplified, e.g., by disregarding decision criteria as seen with *green_deer* (see Sec. V-A). This could stem from fatigue or the already retrieved information being sufficient to make decisions. Users' models could be strongly shaped by the available features of regularly used tools, e.g., some participants examined different quality measures than those described in their vIMM. They might not have found the usually considered one or use whatever measures a tool provides. Our limited time frame also prevented users from entirely conducting their usual process and might have led to adaptations of conducted workflow and applied effort [52]. Instead, users focused on part of their ideal execution, possibly their models' most important parts [53].

In the post-task questionnaire, one participant mentioned trying to replicate their usual process, which led to problems with getting around and getting used to the workflow with SchenQL. Another participant voiced a similar thought in the conduction of T_{ex}; they checked examples for "what can I do with this tool" instead of applying their usual strategy.

Another problem might be the participants' tendency to over-model the vIMMs. One participant answered in the post-task questionnaire "I strongly idealized my search behaviour. (...) my real search behaviour is much simpler".

F. Discussion

1) *Usage of BPMNs*: From our sessions with participants, we observed that the constructed BPMNs were a suitable mean to discuss their task conduction strategies, that these models were generally understood and a clear way to present complex thought processes back to the user (as seen by modifications or clarifications in the verification step of our study). Discussions between the SchenQL experts for constructing the vPGM as well as those between the expert in digital libraries and the SchenQL expert for RQ₅ highlighted the used BPMNs' suitability to formalise the collected data. Our analysis of research questions based on the constructed BPMNs further proves the approach's potential for DL systems' evaluations.

From observations related to the PCMs we found that Law et al.'s [28] method alone has limited suitability for modelling task conduction from think-aloud protocols if users tend to not verbalise crucial actions, even if they are regularly reminded to do so. Additional annotation data, which could be collected from eye-tracking, could help mitigate this problem.

2) *Results of the Evaluation*: When participants described their usual task conduction, they relied on multiple systems and willingly switched tools. Reasons for these switches might

be habit or the tools not offering all required functionalities. When using the example system, users would not be forced to switch tools for a large portion of their usual task conduction, as seen by the considerable part of vIMMs translated to SchenQL. Additionally, many subsequent steps from vIMMs could have been modelled as a single step using SchenQL. This leads to the conclusion that the exemplary evaluated system SchenQL was built for one-shot queries rather than complex processes.

When participants were asked about their usual solutions for the tasks, they described complex processes, possibly caused by the Hawthorne effect [54], but when using a system, they show a "now or never" mentality and only conducted parts of their plan. As possible reasons for this, we assume adjustments made out of time constraints of the evaluation [52, 53, 55], accessibility of features [55], the off-throwing syntax errors, the unknown language/system, missing data and maybe users' information need being satisfied by using a single decision criterion so they did not need to investigate further. A verbalisation of a task conduction could differ from a user's mental conception [56], a BPMN depicting the verbalisation again blurs this model. The model could also change the more a participant finds out about their objective during the search [57]. Generally, vIMMs seemed to be over-modelled.

By the number of problems encountered in the PCMs, participants seemed to have more trouble solving T_{ex} than T_{pa} . This could be attributed to learning effects or task T_{ex} seemingly not to be an everyday information need for part of the participants. In general, the PCMs show users' tendency to iterate over queries multiple times, but only some participants consulted the documentation or example queries for help.

The vIMMs showed participants not solely relying on papers' references but also looking at the context of the references of papers as provided by WoS [34] or scite.

Contrasting previous evaluations of the used example system [6], we did not encounter design fixation [58]. Even though participants were used to their individual workflows, they adapted to the possibilities of the example DL. When asked what users disliked or missed in the current system they did not appear to be stuck to the current design.

Improvements for our example system can function as markers to be incorporated into any DL which strives to better support its users: A system could include a shopping cart style option to save temporary results following Daffodil [11]. Ranking possibilities for results, even soft ranking constraints and diversity aspects should be incorporated. Users should be suggested decompositions of their queries into parallel tasks and follow-up, likely queries to enable simple query specification [59], e.g., recommend a query searching for most cited papers if a user has searched papers.

3) *Cost-Effectiveness*: The monetary cost of the proposed method is non-negligible: Users are required to take part in multiple sessions (in our case, 1.5h per person), which need to be transcribed. Nevertheless, this study type helps uncover discrepancies between users' typical task solution strategies and the functions of a DL with few study participants already.

VII. CONCLUSION

To investigate our general research question *How can we compare users' conceptions of search tasks in digital library with capabilities of such a system?*, we defined five more fine-grained research questions which we observed through a user study with thirteen participants on two everyday tasks in digital libraries: Expert search and relevant paper search. We formalised study participants' usual task conduction processes in BPMNs, the ideal translation of these general processes to an exemplary DL system (SchenQL), and persons actually using the system for a qualitative DL evaluation.

Our experiment found the models to be beneficial and suitable for systematic qualitative DL evaluations. We were able to answer our research questions which investigated different aspects of users' ideal and actual task conduction (RQ₁, RQ₂), limitations of systems (RQ₃), a system's suitability for advanced tasks (RQ₄) and discovered discrepancies between users' perceptions and capabilities of systems (RQ₅). In detail: (1) Users were willing to switch working spheres by explicitly using multiple tools to solve everyday tasks in DLs with their usual exploration strategy. (2) We observed users heavily relying on usage examples for the unknown DL and searching for papers on their chosen topic as entry point for their exploration tasks. (3) DL systems such as SchenQL cannot model interaction with experts, the incorporation of multiple data sources or the type of broad, explorative search as seen with general search engines. (4) Some study participants were able to complete the tasks using the unknown system by translating or adapting their ideal conduction models. (5) Participants tend to overmodel their usual task execution strategy and focus on a portion of the solution process while using SchenQL.

We publish a reusable dataset (see Zenodo [30]) resulting from our evaluation with transcripts of user interviews and the BPMN models, which would be classified as level-4 according to the 5-level system of Gäde et al. [60]. It can be reused to construct new DL systems with requirements from vIMMs or viewing the vIMMs as a qualitative benchmark to investigate systems' compatibility with real user needs. New researchers can apply vIMMs as a starting point for deriving their individual information-seeking strategies. Additionally, one could try to use the BPMNs to construct usage models of different types of users of DLs with the prospect of formulating more real and diverse user types in simulation studies. One could also compare the BPMNs to existing search stratagems [46, 61].

Future work could apply our proposed qualitative evaluation technique to other DL search systems and combine these results in a meta-study. For prospective studies, eye-tracking software could be used to try to better explain user behaviour and annotate the process models. Different formal representations of user models could be incorporated and compared.

When focusing more on the results of our exemplary evaluation, users' readiness to switch systems could be investigated. Investigating potential links between study participants' expertise with DLs or their chosen topics and their task conduction could shed light on their individual cognitive loads [62].

REFERENCES

- [1] C. Betts, J. Power, and W. Ammar, "Grapal: Connecting the dots in scientific literature," in *ACL '19*. ACL, 2019, pp. 147–152. [Online]. Available: <https://doi.org/10.18653/v1/p19-3025>
- [2] H. Kroll, F. Plötzky, J. Pirklbauer, and W. Balke, "What a publication tells you: benefits of narrative information access in digital libraries," in *JCDL '22*. ACM, 2022, p. 9. [Online]. Available: <https://doi.org/10.1145/3529372.3530928>
- [3] C. K. Kreutz, M. Blum, and R. Schenkel, "Schenql: a query language for bibliographic data with aggregations and domain-specific functions," in *JCDL '22*. ACM, 2022, pp. 37:1–37:5. [Online]. Available: <https://doi.org/10.1145/3529372.3533282>
- [4] I. Xie, S. Joo, and K. Matusiak, "Digital library evaluation measures in academic settings: Perspectives from scholars and practitioners," *Journal of Librarianship and Information Science*, vol. 53, p. 096100062093550, 06 2020.
- [5] C. C. Kuhlthau, *Seeking Meaning: A Process Approach to Library and Information Services*, ser. Information management, policy, and services. Libraries Unlimited, 2004. [Online]. Available: <https://books.google.de/books?id=feDgAAAAMAAJ>
- [6] C. K. Kreutz, M. Wolz, J. Knack, B. Weyers, and R. Schenkel, "Schenql: in-depth analysis of a query language for bibliographic metadata," *Int. J. Digit. Libr.*, vol. 23, no. 2, pp. 113–132, 2022. [Online]. Available: <https://doi.org/10.1007/s00799-021-00317-8>
- [7] Ö. Dalkıran, İ. Aker, S. Öztemiz, Z. Taskin, and S. Tunç, "Usability testing of digital libraries: The experience of eprints," *Procedia - Social and Behavioral Sciences*, vol. 147, 08 2014.
- [8] Y. Zhu, M. C. Kim, and E. Yan, "Evaluating interactive bibliographic information retrieval systems: A user-centered approach," *Proceedings of the Association for Information Science and Technology*, vol. 55, pp. 628–637, 01 2018.
- [9] J. Dinet, M. Favart, and J.-M. Passerault, "Searching for information in an online public access catalogue (opac): the impacts of information search expertise on the use of boolean operators," *J. Comput. Assist. Learn.*, vol. 20, pp. 338–346, 2004.
- [10] S. W. Kumpulainen and H. Kautonen, "Accidentally successful searching: Users' perceptions of a digital library," in *CHIIR '17*. ACM, 2017, p. 257–260. [Online]. Available: <https://doi.org/10.1145/3020165.3022124>
- [11] A. Schaefer, M. Jordan, C. Klas, and N. Fuhr, "Active support for query formulation in virtual digital libraries: A case study with DAFFODIL," in *ECDL '05*, ser. LNCS, vol. 3652. Springer, 2005, pp. 414–425. [Online]. Available: https://doi.org/10.1007/11551362_37
- [12] T. Krämer, A. Papenmeier, Z. Carevic, D. Kern, and B. Mathiak, "Data-seeking behaviour in the social sciences," *Int. J. Digit. Libr.*, vol. 22, no. 2, pp. 175–195, 2021. [Online]. Available: <https://doi.org/10.1007/s00799-021-00303-0>
- [13] W. Fernandes and B. Cendon, "A study of non-users of digital libraries: the case of the capes digital library in brazil," *The Electronic Library*, vol. ahead-of-print, 07 2021.
- [14] M. Miller, G. Choi, and L. Chell, "Comparison of three digital library interfaces: Open library, google books, and hathi trust," in *JCDL '12*. ACM, 2012, p. 367–368. [Online]. Available: <https://doi.org/10.1145/2232817.2232894>
- [15] S. Liang, D. He, D. Wu, and H. Hu, "Challenges and opportunities of acm digital library: A preliminary survey on different users," in *Sustainable Digital Communities*. Springer, 2020, pp. 278–287.
- [16] V. Bartalesi, C. Meghini, D. Metilli, and P. Andriani, "Usability evaluation of the digital library dantesources," in *Human-Computer Interaction. Novel User Experiences*. Springer, 2016, pp. 191–203.
- [17] G. Marchionini, C. Plaisant, and A. Komlodi, "Interfaces and tools for the library of congress national digital library program," *Information Processing & Management*, vol. 34, no. 5, pp. 535–555, 1998. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S030645739800020X>
- [18] D. Kelly, "Methods for evaluating interactive information retrieval systems with users," *Found. Trends Inf. Retr.*, vol. 3, no. 1-2, pp. 1–224, 2009. [Online]. Available: <https://doi.org/10.1561/1500000012>
- [19] M. Kellar, C. Watters, and M. Shepherd, "A field study characterizing web-based information-seeking tasks," *J. Am. Soc. Inf. Sci. Technol.*, vol. 58, no. 7, p. 999–1018, may 2007.
- [20] N. Ford, D. Miller, and N. Moss, "The role of individual differences in internet searching: An empirical study," *Journal of the American Society for Information Science and Technology*, vol. 52, no. 12, pp. 1049–1066, 2001. [Online]. Available: <https://onlinelibrary.wiley.com/doi/abs/10.1002/asi.1165>
- [21] K.-S. Kim and B. Allen, "Cognitive and task influences on web searching behavior," *Journal of the American Society for Information Science and Technology*, vol. 53, no. 2, pp. 109–119, 2002. [Online]. Available: <https://onlinelibrary.wiley.com/doi/abs/10.1002/asi.10014>
- [22] K. Byström, "Information and information sources in tasks of varying complexity," *Journal of the American Society for Information Science and Technology*, vol. 53, no. 7, pp. 581–591, 2002. [Online]. Available: <https://onlinelibrary.wiley.com/doi/abs/10.1002/asi.10064>
- [23] Y. Rogers, H. Sharp, and J. Preece, "Interaction design," *Interaction Design: Beyond Human-Computer Interaction*, Wiley, pp. 1–34, 2011.
- [24] J. Bowen, A. Dittmar, and B. Weyers, "Task modelling for interactive system design: A survey of historical

- trends, gaps and future needs,” *Proceedings of the ACM on Human-Computer Interaction*, vol. 5, no. EICS, pp. 1–22, 2021.
- [25] H. Kautonen, “Playing dice with a digital library: Analysis of an artist using a new information resource for her art production,” in *Human-Computer Interaction: Users and Contexts*. Springer, 2015, pp. 430–440.
- [26] M. T. d. V. Dias, S. W. M. Siqueira, B. Pereira Nunes, M. Bortoluzzi, and I. Marenzi, “Modeling exploratory search as a knowledge-intensive process,” in *ICALT '18*, 2018, pp. 34–38.
- [27] P. Vakkari, “Exploratory searching as conceptual exploration,” 01 2010.
- [28] Y. C. Law, W. Wehrt, S. Sonnentag, and B. Weyers, “Obtaining Semi-Formal Models from Qualitative Data: From Interviews Into BPMN Models in User-Centered Design Processes,” *International Journal of Human-Computer Interaction*, vol. 39, no. 3, pp. 476–493, Feb. 2023. [Online]. Available: <https://www.tandfonline.com/doi/full/10.1080/10447318.2022.2041899>
- [29] E. J. Kelly, “Assessment of digitized library and archives materials: A literature review,” *Journal of Web Librarianship*, vol. 8, no. 4, pp. 384–403, 2014. [Online]. Available: <https://doi.org/10.1080/19322909.2014.954740>
- [30] C. K. Kreutz, S. Myshkina, M. Blum, P. Schaer, R. Schenkel, and B. Weyers, “Formalised information needs dataset,” 2023. [Online]. Available: <https://doi.org/10.5281/zenodo.7826530>
- [31] A. Hotho, R. Jäschke, D. Benz, M. Grahl, B. Krause, C. Schmitz, and G. Stumme, “Social bookmarking am beispiel bibsonomy,” in *Social Semantic Web: Web 2.0 - Was nun?*, ser. X.media.press. Springer, 2009, pp. 363–391. [Online]. Available: https://doi.org/10.1007/978-3-540-72216-8_18
- [32] M. Ley, “DBLP - some lessons learned,” *Proc. VLDB Endow.*, vol. 2, no. 2, pp. 1493–1500, 2009. [Online]. Available: <http://www.vldb.org/pvldb/vol2/vldb09-98.pdf>
- [33] W. Ammar, D. Groeneveld, C. Bhagavatula, I. Beltagy, M. Crawford, D. Downey, J. Dunkelberger, A. Elgohary, S. Feldman, V. Ha, R. Kinney, S. Kohlmeier, K. Lo, T. Murray, H. Ooi, M. E. Peters, J. Power, S. Skjonsberg, L. L. Wang, C. Wilhelm, Z. Yuan, M. van Zuylen, and O. Etzioni, “Construction of the literature graph in semantic scholar,” in *NAACL-HLT (3)*. ACL, 2018, pp. 84–91.
- [34] L. Salisbury, “Web of science and scopus: A comparative review of content and searching capabilities,” *The Charleston Advisor*, vol. 11, 01 2009.
- [35] C. Herzog, D. W. Hook, and S. R. Konkiel, “Dimensions: Bringing down barriers between scientometricians and data,” *Quant. Sci. Stud.*, vol. 1, no. 1, pp. 387–395, 2020. [Online]. Available: https://doi.org/10.1162/qss_a_00020
- [36] J. Priem, H. A. Piwowar, and R. Orr, “Openalex: A fully-open index of scholarly works, authors, venues, institutions, and concepts,” *CoRR*, vol. abs/2205.01833, 2022. [Online]. Available: <https://doi.org/10.48550/arXiv.2205.01833>
- [37] F. Paternò, “Concurtasktrees: An engineered approach to model-based design of interactive systems,” 07 2008.
- [38] C. Martinie, P. Palanque, and M. Winckler, “Structuring and Composition Mechanisms to Address Scalability Issues in Task Models,” in *INTERACT '11*, ser. LNCS, vol. 6948, no. Part 3. Springer, Sep. 2011, pp. 589–609. [Online]. Available: <https://hal.inria.fr/hal-01591816>
- [39] M. T. d. V. Dias, “Understanding web search patterns through exploratory search as a knowledge-intensive process.” 2019. [Online]. Available: [http://www.repositorio-bc.unirio.br:8080/xmlui/bitstream/handle/unirio/12894/DIAS%20-%20MARCELO%20TIBAU%20DE%20VASCONCELLOS%20\(1\).pdf](http://www.repositorio-bc.unirio.br:8080/xmlui/bitstream/handle/unirio/12894/DIAS%20-%20MARCELO%20TIBAU%20DE%20VASCONCELLOS%20(1).pdf)
- [40] T. Allweyer, *BPMN 2.0 - Business Process Model and Notation: Einführung in den Standard für die Geschäftsprozessmodellierung*. Books on Demand, 2020.
- [41] A. Soufan, I. Ruthven, and L. Azzopardi, “Searching the literature: An analysis of an exploratory search task,” in *CHIIR '22*. ACM, 2022, pp. 146–157. [Online]. Available: <https://doi.org/10.1145/3498366.3505818>
- [42] A. Sinha, Z. Shen, Y. Song, H. Ma, D. Eide, B. P. Hsu, and K. Wang, “An overview of microsoft academic service (MAS) and applications,” in *WWW '15 (Companion Volume)*. ACM, 2015, pp. 243–246.
- [43] J. Tang, J. Zhang, L. Yao, J. Li, L. Zhang, and Z. Su, “Arnetminer: extraction and mining of academic social networks,” in *KDD '08*. ACM, 2008, pp. 990–998.
- [44] F. Zhang, X. Liu, J. Tang, Y. Dong, P. Yao, J. Zhang, X. Gu, Y. Wang, B. Shao, R. Li, and K. Wang, “OAG: toward linking large-scale heterogeneous entity graphs,” in *KDD '19*. ACM, 2019, pp. 2585–2595.
- [45] T. Dresing and T. Pehl, *Praxisbuch Interview, Transkription & Analyse Anleitungen und Regelsysteme für qualitativ Forschende*, 2012. [Online]. Available: <https://books.google.de/books?id=tzQIR8HI36QC>
- [46] M. J. Bates, “The design of browsing and berrypicking techniques for the online search interface,” *Online Review*, vol. 13(5), pp. 407–424, 1989.
- [47] T. Khazaei and O. Hoerber, “Metadata visualization of scholarly search results: supporting exploration and discovery,” in *I-KNOW '12*. ACM, 2012, p. 21. [Online]. Available: <https://doi.org/10.1145/2362456.2362483>
- [48] G. Mark, V. M. Gonzalez, and J. Harris, “No task left behind? examining the nature of fragmented work,” in *CHI '05*. ACM, 2005, pp. 321–330.
- [49] S. Jeuris and J. E. Bardram, “Dedicated workspaces: Faster resumption times and reduced cognitive load in sequential multitasking,” *Computers in Human Behavior*, vol. 62, pp. 404–414, 2016. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0747563216302308>
- [50] X. Yuan and N. J. Belkin, “Supporting multiple

- information-seeking strategies in a single system framework,” in *SIGIR '07*. ACM, 2007, pp. 247–254. [Online]. Available: <https://doi.org/10.1145/1277741.1277786>
- [51] P. Vakkari and N. Hakala, “Changes in relevance criteria and problem stages in task performance,” *J. Documentation*, vol. 56, no. 5, pp. 540–562, 2000. [Online]. Available: <https://doi.org/10.1108/EUM000000007127>
- [52] A. Crescenzi, R. G. Capra, B. Choi, and Y. Li, “Adaptation in information search and decision-making under time constraints,” in *CHIIR '21*. ACM, 2021, pp. 95–105. [Online]. Available: <https://doi.org/10.1145/3406522.3446030>
- [53] A. Edland and O. Svenson, *Judgment and Decision Making Under Time Pressure*, 01 1993, pp. 27–40.
- [54] C. R. Harrell, B. Gladwin, and M. P. Hoag, “Mitigating the “hawthorne effect” in simulation studies,” in *WSC '13*, 2013, pp. 2722–2729.
- [55] A. Crescenzi, R. Capra, and J. Arguello, “Time pressure, user satisfaction and task difficulty,” in *ASIS&T '13*, ser. Proc. Assoc. Inf. Sci. Technol., vol. 50, no. 1. Wiley, 2013, pp. 1–4. [Online]. Available: <https://doi.org/10.1002/meet.14505001121>
- [56] R. S. Taylor, “Question-negotiation and information seeking in libraries,” *Coll. Res. Libr.*, vol. 76, no. 3, pp. 251–267, 2015. [Online]. Available: <https://doi.org/10.5860/crl.76.3.251>
- [57] C. Cole, “A theory of information need for information retrieval that connects information to knowledge,” *Journal of the American Society for Information Science and Technology*, vol. 62, no. 7, pp. 1216–1231, 2011. [Online]. Available: <https://onlinelibrary.wiley.com/doi/abs/10.1002/asi.21541>
- [58] D. G. Jansson and S. M. Smith, “Design fixation,” *Design Studies*, vol. 12, no. 1, pp. 3–11, 1991. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/0142694X9190003F>
- [59] J. Liu, C. Liu, X. Yuan, and N. J. Belkin, “Understanding searchers’ perception of task difficulty: Relationships with task type,” in *ASIS&T '11*, ser. Proc. Assoc. Inf. Sci. Technol., vol. 48, no. 1. Wiley, 2011, pp. 1–10. [Online]. Available: <https://doi.org/10.1002/meet.2011.14504801152>
- [60] M. Gäde, M. Koolen, M. Hall, T. Bogers, and V. Petras, “A Manifesto on Resource Re-Use in Interactive Information Retrieval,” in *CHIIR '21*. ACM, Mar. 2021, pp. 141–149. [Online]. Available: <https://dl.acm.org/doi/10.1145/3406522.3446056>
- [61] A. Kacem and P. Mayr, “Analysis of search stratagem utilisation,” *Scientometrics*, vol. 116, no. 2, pp. 1383–1400, 2018. [Online]. Available: <https://doi.org/10.1007/s11192-018-2821-8>
- [62] J. Sweller, “Cognitive load during problem solving: Effects on learning,” *Cogn. Sci.*, vol. 12, no. 2, pp. 257–285, 1988. [Online]. Available: https://doi.org/10.1207/s15516709cog1202_4