# Studying Databases of Intentions: Do Search Query Logs Capture Knowledge about Common Human Goals?

**Markus Strohmaier**
Graz University of Technology and Know-Center
Graz, Austria
markus.strohmaier@tugraz.at

**Mark Kröll**
Graz University of Technology
Graz, Austria
mkroell@tugraz.at

## ABSTRACT

Access to knowledge about common human goals has been found critical for realizing the vision of intelligent agents acting upon user intent on the web. Yet, the acquisition of knowledge about common human goals represents a major challenge. In a departure from existing approaches, this paper investigates a novel resource for knowledge acquisition: The utilization of search query logs for this task. By relating goals contained in search query logs with goals contained in existing commonsense knowledge bases such as ConceptNet, we aim to shed light on the usefulness of search query logs for capturing knowledge about common human goals. The main contribution of this paper consists of insights generated from an empirical study comparing common human goals contained in two large search query logs (AOL and Microsoft Research) with goals contained in the commonsense knowledge base ConceptNet. The paper sketches ways how goals from search query logs could be used to address the goal acquisition and goal coverage problem related to commonsense knowledge bases.

## Categories and Subject Descriptors

H.3.1 [**Information Storage and Retrieval**]: Content Analysis and Indexing; I.2.6 [**Artificial Intelligence**]: Learning; I.2.7 [**Artificial Intelligence**]: Natural Language Processing

## General Terms

Algorithms, Experimentation

## Keywords

Knowledge acquisition, human goals, commonsense knowledge

## INTRODUCTION

To realize the vision of common-sense enabled and goal-oriented agents on the web, agents must have programmatic access to the set and variety of common human goals, in order to reason about them and to provide services that help satisfy users' needs [14][23]. In Berner-Lee's vision, an agent aiming to, for example, "plan a trip to Vienna" would need to have some means to understand that "plan a trip" is likely to involve a set of other goals or services, such as "contact a travel agency" and "book a hotel". This type of knowledge has been characterized as commonsense knowledge, i.e. knowledge that humans are generally assumed to possess, but which is extremely difficult for computers to acquire [16]. Having such knowledge available is a prerequisite for applications such as GOOSE – a commonsense-enabled search engine [15] - or EventMinder – a commonsense-enabled calendar application [25], which are inspiring prototypes for the potential of traditional applications augmented with commonsense knowledge.

Current research projects aiming to capture and organize commonsense knowledge include CyC [11] or Openmind / ConceptNet [24]. These projects utilize human knowledge engineering [11], volunteer-based [16], game-based [14] or semi-automatic approaches [5] for knowledge acquisition, where common human goals can be considered a subset of the enormous breadth of commonsense knowledge. However, existing attempts to capture knowledge about common human goals generally suffer from a number of problems, including: 1) the *goal acquisition problem (or bottleneck)*, which refers to the costs associated with knowledge acquisition [14] and 2) the *goal coverage problem*, which refers to the difficulty of capturing the tremendous variety and range in the set of common human goals [5]. These problems have hindered progress in capturing broad knowledge about common human goals, and have hindered the development of intelligent agents and applications on the web. In this paper, we are seeking to explore the utility of search queries, i.e. user needs expressed in small textual fragments, to help address the above challenges.

Using textual contributions of volunteers on the web for knowledge acquisition purposes is not a new idea (cf. MIT's Openmind initiative [16]). However, in contrast to volunteer-based systems, Search Query Logs provide an abundant and seemingly endless stream of information about human needs. This has led query logs being referred to as "Databases of Intentions" in the past. Databases of intentions[1] refer to the observation that Search Query Logs provide unique, up-to-date and detailed insights into human motivations, needs and goals (such as "buying a house") that can be analyzed, studied and used for different purposes. While search query logs have been utilized successfully for knowledge acquisition in a range of different contexts [18], they have not been used to capture explicit knowledge about human goals, partly because query logs pose a num-

---

[1] http://battellemedia.com/archives/000063.php, last accessed on April 15, 2009

ber of challenges: The majority of search queries is usually short [18] and ambiguous [1], they convey user goals at different degrees of intentional explicitness [27], and they are often consisting of arbitrary concatenations of terms that frequently contain misspellings. Yet, recent research revealed that a number of search queries actually contain explicit statements of human goals [27], and that the space of queries in search query logs is vast and topically diverse [19]. This would make query logs a seemingly attractive resource for acquiring knowledge about a diverse range of common human goals. *But how useful are Search Query Logs for capturing common human goals?* In this paper, we are investigating the following guiding research questions related to this task:

**RQ1**: Do Search Query Logs contain knowledge about common human goals?
**RQ2**: If they do, what is the nature of common human goals shared by ConceptNet and two large Query Logs?
**RQ3**: Do goals contained in ConceptNet and Search Query Logs differ w.r.t. scope, and if so - how?
**RQ4**: Can goals contained in Search Query Logs be used to refine ConceptNet?
**RQ5**: Can goals contained in Search Query Logs be used to add novel nodes and relations to ConceptNet?

Finding answers to these questions would help assess the usefulness of search query logs to lower the costs often associated with commonsense knowledge capture, and would help to assess the potential to improve the coverage of existing knowledge bases such as ConceptNet.

## Contributions

In previous work, we have developed an automatic classification approach focused on identifying search queries that contain explicit statements of human goals [26]. In this paper, we apply our approach to two very large query logs provided by AOL and Microsoft Research, yielding a set of ~115.000 queries that contain common human goals in 77 out of 100 cases (77% precision).

The paper makes the following contributions: Using the set of goals acquired from search query logs, we conduct comparative analyses to assess the nature of these goals in relation to an exemplary commonsense knowledge base, in our case ConceptNet. Subsequently, we discuss preliminary proposals on how to use common human goals acquired from query logs to refine and expand existing knowledge bases. Our findings suggest that goals acquired from search query logs have the potential to expand coverage of common human goals in existing knowledge bases while maintaining reasonable precision scores. To the best of our knowledge, this work represents the first *comparative* study of the utility of search query logs for capturing common human goals. Our findings suggest that search query logs indeed represent a *viable*, yet largely *untapped* alternative for acquiring knowledge about common human goals.

## AUTOMATIC GOAL ACQUISITION FROM SEARCH QUERY LOGS

In this paper, we use a machine learning approach previously developed by the authors of this paper [26] for the purpose of automatically classifying queries regarding whether they contain an explicit goal statement or not. This approach can be used to acquire a set of queries containing goals from search query logs. Table 1 gives some examples of actual queries containing/not containing statements of human goals (obtained from [19]).

**Table 1 Examples of queries obtained from [19]**

| Queries containing explicit goal statements | Queries not containing explicit goal statements |
|---|---|
| "sell my car" | "Mazda dealership" |
| "play online poker" | "online games" |
| "find home to rent in Florida" | "Miami beach houses" |
| "passing a drug test" | "drug test" |
| "raising your credit score" | "credit cards" |

Based on work that highlights the crucial role of verbs in explicit statements of goals ([13],[22]), we define queries containing human goals in the following way:

*A search query is regarded to contain an explicit user goal whenever the query 1) contains at least one verb and 2) describes a plausible state of affairs that the user may want to achieve or avoid (cf. [22]) in 3) a recognizable way [26].*

"Recognizable" refers to what [9] defines as "trivial to identify" by a subject within a given attention span. "Plausible" refers to an external observer's assessment whether the goal contained in a query could likely represent the goal of a user who formulates the given query. This definition has been shown to produce reasonable inter-rater agreements among independent subjects [26]. It is important to note that it would be rather difficult to assess queries solely based on data from an anonymous query log due to the inherent *goal verification problem* of such a task [26]. However, the objectives of this work are more modest: In this paper, we are interested in acquiring *plausible* common human goals for *knowledge capture* purposes. An advantage of acquiring broad knowledge about *plausible* goals is that it can put constraints on the space of *all* goals, which plays a role in, for example, goal recognition [6] or query disambiguation [1].

"Queries containing explicit goals" according to our definition thus can be related to what other researchers have characterized as "better queries", or queries that have "more precise goals" (R. Baeza-Yates at the "Future of Web Search" Workshop 2006, Barcelona). A query does not contain an explicit goal when it is difficult or extremely hard to elicit some specific goal from the query. Examples include blank queries, or queries such as "car" or "travel", which embody goals on a very general, ambiguous and mostly implicit level.

## STUDY SETUP

Our comparative study is based on analyzing two different types of datasets, i.e. Search Query Logs and commonsense knowledge bases.

## Description of Method: A Comparative Study

In order to assess the degree to which search query logs contain knowledge about common human goals, we pursue a comparative study approach. Our approach aims to investigate the types of common human goals contained in search query logs *in relation to* existing commonsense knowledge bases. For that, we extract a large number of common human goals from search query logs and commonsense knowledge bases and investigate their characteristics, their differences and commonalities. Subsequently, we are interested in exploring the ways of interlinking goals acquired from search query logs with commonsense knowledge bases, as this might represent a way of overcoming some of the challenges introduced earlier. However, given the state of research, it would be premature to propose definite ways of expanding existing knowledge bases with goals acquired from query logs. Instead, this work follows an exploratory research approach, investigating the range of possibilities of relating knowledge about common human goals from query logs and commonsense knowledge bases, and assessing their usefulness.

## Description of Datasets

**Search Query Logs**: We applied our automatic goal acquisition method to two large search query logs recorded by AOL and Microsoft research in 2006. We combined the data from two query logs to (i) increase the number of samples as well as (ii) to avoid a potential bias that is introduced by using only one search query log. The first query log, the MSN search query log excerpt, contains about 15 million queries (from US users) that were sampled over one month in May, 2006. The second log, the AOL search query log [19], contains ~ 20 million search queries collected from 657,426 unique user ID's between March 1, 2006 and May 31, 2006. Search queries from both logs were extracted using the same method, and underwent various normalization steps, e.g. trimming of each query and space sequence reduction to one space character. A typical entry in a search query log contains the following items: {UserID, Query, Timestamp, (ItemRank, URL)*}. From these two search query logs combined, a set of ~115.000 human goals was extracted by our automatic classification method. The resulting set of queries containing goals has a precision of 77% [26] – which is comparable to the precision of commonsense knowledge bases [5].

**ConceptNet**: To enable comparative analysis, we chose ConceptNet 2.1 [16] as an example of a commonsense knowledge base, mainly due to its open availability and its informal knowledge representation, which makes it particularly suitable for our purposes. Knowledge in ConceptNet is partly represented in free-form text facilitating the com-

parison with search queries. In order to compare goals from Search Query Logs with goals in ConceptNet, we had to compile a subset of entries from ConceptNet that focus on goals. For that, we considered entries that are connected by ConceptNet's "MotivatedByGoal", "UsedFor" and "CapableOf" relations. We imposed the following restrictions on all phrases: Goal phrases had to contain at least one verb and at least one object. The restriction was enforced by examining the POS tags[2] of the phrases. This resulted in a set of ~68.000 goals acquired from ConceptNet.

The AOL dataset was obtained from a secondary resource[3] while the MS dataset was provided by Microsoft Research.

## RESULTS

In the following, we present our findings related to the five research questions introduced earlier.

## Existence of Common Human Goals

> **RQ 1: Do Search Query Logs contain knowledge about common human goals?**
>
> **Result**: Search Query Logs share a significant number of common human goals with goals contained in ConceptNet.
>
> **Conclusion**: Search Query Logs are a potential source of common human goals.

To answer research question 1, we calculated an intersection set between the search query log and ConceptNet goal datasets. An adequate number of common entries would indicate the presence of commonsense knowledge in search query logs. To calculate intersection, we devised a simple goal matching algorithm (similar to [14]) to identify matching pairs of human goals from the two different corpora. All entries were pre-processed, i.e. stop words were removed and all remaining tokens were stemmed[4]. This is illustrated in the examples below:

| Before Pre-Processing | After Pre-Processing |
|---|---|
| to protect your family | ['protect', 'famili'] |
| have something to do during breakfast | ['have', 'do', 'breakfast'] |
| how to tell kids about suicide | ['tell', 'kid', 'suicid'] |

For two goals to match, they had to contain an equal numbers of identical stems. Our algorithm focused on lexical characteristics only. Table 2 illustrates some examples of matching and non-matching entries from ConceptNet and search query logs.

We decided to pursue a traditional bag-of-words approach where sequence information is neglected. The implications of this approach are discussed in the limitations section of this paper. The algorithm yielded ~2300 ConceptNet goals and ~3100 search query log goals (occurrences) that pro-

---

[2] Stanford Part-Of-Speech Tagger version 1.6

[3] http://www.gregsadetsky.com/aol-data/, from July 15th, 2007

[4] NLTK Porter Stemmer

duced positive matches. While the number of common goals appears small, our findings are based on comparably small samples of query logs and nevertheless provide first evidence that search query logs contain commonsense knowledge to some extent.

**Table 2: Matching and Non-Matching ConceptNet and Search Query Log Entries. The matching sequences are highlighted.**

| ConceptNet Goals | Search Query Log Goals | Match |
|---|---|---|
| make paper airplanes | how to make paper airplanes | yes |
| get in shape | getting into shape | yes |
| we buy houses | buy a house | yes |
| we buy houses | purchase a house | no |
| make money | make more money | yes |
| make money | make online money | no |

Next we are interested in studying the nature of goals contained in both search query logs and ConceptNet.

## Overlap Between Datasets

**RQ 2: What is the nature of goals shared by Concept-Net and two large Query Logs?**

**Result**: "Build", "Obtain" and "Perform" goals are types of goals frequently found in both corpora. Other types of goals are less shared.

**Conclusion**: Overlap between goals in query logs and ConceptNet focuses on a few prominent types of goals.

In order to learn more about common human goals shared by the two different corpora, we categorized the ~3100 search query log goals into a subset of Levin's verb class taxonomy [12]. We selected verb classes that we deemed relevant for reflecting human activities and discarded those with few entries afterwards such as 'complain' (Verb class "37.8"). Figure 1 shows the verb class histogram.
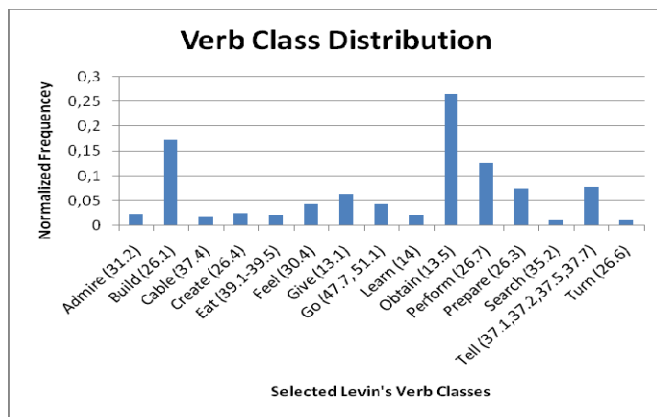


**Figure 1: Distribution over Selected Levin's Verb Classes of Goals in the Intersection Set. Corresponding Levin's verb classes are denoted in brackets.**

The three dominant verb classes 'build', 'obtain' and 'perform', can easily be identified. Their dominance can be explained by frequent occurrences of verbs such as 'make' (class 'build'), 'get, find, and buy' (class 'obtain') and 'take and play' (class 'perform'). These verbs are part of common human goals such as 'make money', 'buy food' or 'play an instrument'. Another outstanding verb class is 'tell'. We expected this class to be prominent for containing verbs that describe various ways of communication such as 'explain', 'narrate' or 'suggest'. Overall, the distribution of overlapping types of goals suggests that search query logs do not equally cover different types of goals contained in ConceptNet. This shows that search query logs differ with regard to the scope of common human goals expressed in them. To further study this question, we investigated differences between the two datasets in greater detail.

## Differences between Datasets

**RQ 3: Do goals contained in ConceptNet and Search Query Logs differ w.r.t. scope, and if so - how?**

**Result:** Goals contained in ConceptNet and Search Query Logs exhibit similar scope, but exhibit differences with regard to the extent of coverage.

**Conclusion:** Search Query Logs might represent a useful resource to expand existing knowledge bases such as ConceptNet.

As in the previous section, we generated a verb class histogram for the complementary sets of goals, i.e. the set of goals in both datasets that is not contained in the intersection set. Our initial intuition was that ConceptNet goals would be biased towards everyday situations and human characteristics such as *eating*, *feeling* and *living*. This intuition was confirmed by our results.
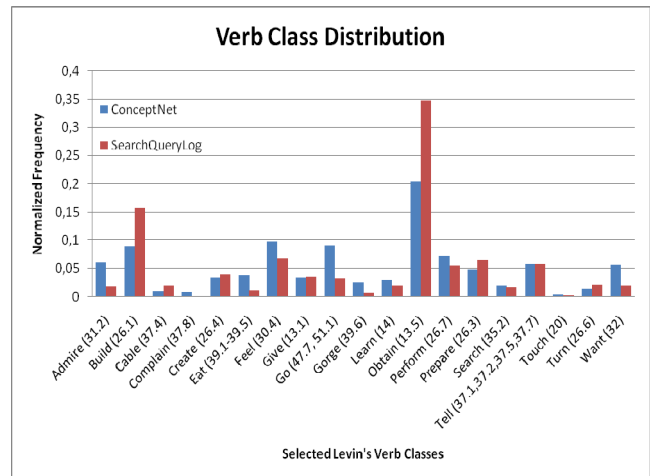


**Figure 2: Distribution over Selected Levin's Verb Classes of Goals in the two Complementary Sets. Corresponding Levin's verb classes are denoted in brackets.**

The verb histogram in Figure 2 shows that verb classes 'eat', 'gorge', 'touch', and 'feel' are more prominent in the ConceptNet set. Similarly, we expected the classes 'search' and 'obtain' to be the dominating search query log goals. This can only be partly observed in our results: Levin's

verb class 'obtain' dominates the goals acquired from search query logs, which contain frequently occurring verb instances such as 'get', 'buy' and 'find'. However, the verb class 'search' is not dominant in our dataset. We believe this class is underrepresented due to the fact that search engines already represent a means for searching the web. As a consequence, search queries do not seem to contain explicit verbs describing search intent itself.

The verb class histogram in Figure 2 reveals that each dataset favors a different range of human activities. In fact, this suggests that query logs could actually *help increase coverage of commonsense knowledge bases*, for example by focusing on types of goals that are more prevalent in search query logs, such as 'obtain' and 'build'. In the following subsection, we are interested in this question. We study whether and how we could tap into search query logs for refining and expanding commonsense knowledge bases with common human goals contained in search query logs. We conducted two investigations to explore how goals acquired from Query Logs could be integrated with ConceptNet.

## Refining ConceptNet Using Query Logs

RQ 4: Can goals contained in Search Query Logs be used to refine ConceptNet?

**Preliminary Result**: By example, we have shown that simple, but effective, refinements can be done.

**Conclusion**: Search query logs can add reasonable refinements to commonsense knowledge bases such as ConceptNet.

Because ConceptNet has a general scope, and search queries are often driven by specific needs, it seems reasonable to assume that given a ConceptNet goal, potential goal refinements from search query logs can be retrieved. Refinements could be integrated into ConceptNet by connecting the ConceptNet goal and the refinement with a dedicated relationship introduced for this purpose. The advantage of utilizing search query logs for this task is that it offers frequency information, i.e. how many users issued a certain query? Frequent queries can be expected to be of reasonable content, thereby minimizing the risk of selecting inappropriate candidates. This is an idea similar to the 'f-value', an internal parameter in ConceptNet that is intended to be related to the quality of a [node, relation, node] triple.

In the following, we envision a simple potential approach that aims at retrieving refinement candidates from search query logs. In the proposed approach, search queries are pre-processed as described previously. Prospective refinement candidates are required to contain exactly the same tokens (stems) as the corresponding ConceptNet goal. In addition, we can require candidates to contain one adjective that must not be at the beginning or the end of the phrase. By applying this strict but simple search pattern, we expected to retrieve suitable refinement candidates of high quality. To investigate this idea, we extracted candidate refinements for all ConceptNet goals in the intersection set

yielding an average of ~3 refinements per ConceptNet goal. Note that this rate would represent a *significant quantitative, and relevant refinement of the goals contained in ConceptNet*. Table 3 illustrates exemplary candidate refinements and their corresponding goals. The attached frequency information is related to the occurrences of this query in the search quer y logs.

**Table 3: Exemplary Refinement Candidates for Selected ConceptNet Goals.**

| ConceptNet Goal | List of Refinements From Query Logs |
| --- | --- |
| now buy this car | buy new car (20), buy a cheap car (2), buying rental cars (2), buy electric car (2), buy a used car (19), buy old cars (3), buying wise car (1) |
| finding friends | to find old friends (4), find high school friends(1), find lost friend (11), find best friends (1), find elementary school friend (1), find free online friends (1), find past military friends (1) |
| writing a paper | write an argumentative paper (1), write an informative paper (1), write an autobiographical paper (1), write a narrative paper (1) |
| cutting your hair | cutting my own hair (1), cut short hair (3), cut black hair (1), cutting long hair (1) |
| feeding the baby | feeding a newborn baby (1) |
| find a partner | finding sexual partners (1) |
| train a dog | train an abused dog (1), train a deaf dog (1) |
| making coffee | make perfect coffee (1), make a flaming coffee (1) |

Using query frequency information, frequency information could act as a filter for identifying the most suitable candidates for refinements. As illustrated in Table 3, higher frequencies could indicate candidates that are better suited to refine ConceptNet goals than low frequency refinements. Thereby, frequency information can be used to rank candidate refinements und thereby exclude potentially inappropriate refinements. Table 4 compares a small sample of high frequency candidates to low frequency candidates.

**Table 4: High and Low Frequency Candidates are contrasted. Frequency information seems to be a useful quality parameter to assess commonsense knowledge.**

| ConceptNet Goal | High Frequency Candidate | Low Frequency Candidate |
| --- | --- | --- |
| now buy this car | buy a new car (20), buy a used car (19) | buying wise car (1) |
| finding friends | find lost friend (11), find old friends (4) | find past military friends (1), find free online friends (1) |

Although we have not investigated the issue of finding suitable thresholds, it seems that frequency information can serve as a proxy indicator for automatically ranking and reviewing potential ConceptNet refinements from human goals contained in Search Query Log. Next, we are interested whether goals acquired from search query logs can add new and interesting information to commonsense knowledge bases.

## Extending ConceptNet Using Query Logs

> **RQ 5: Can goals contained in Search Query Logs be used to add novel nodes and relations to ConceptNet?**
>
> **Result**: ConceptNet can be expanded with novel goals acquired from search query logs at a precision that is comparable to the precision reported for ConceptNet.
>
> **Conclusion**: Search query logs can add useful knowledge that is not yet contained in existing knowledge bases.

Amongst ConceptNet's numerous relation types, the '*MotivatedByGoal*' relation establishes an intentional connection between two nodes in ConceptNet. ['read a book', 'MotivatedByGoal', 'learn something'] or ['wait tables', 'MotivatedByGoal', 'make money'] are two exemplary ConceptNet instances of this particular relation type. These triples typically comprise a 'LeftText', which describes an action, and a 'RightText', which describes the motivation behind the action.

Building on this, we are interested in identifying whether goals from Search Query Logs could represent proper motivations (RightTexts) for ConceptNet actions (LeftTexts). To investigate this question, we started from 'MotivatedByGoal' triples that were already contained in ConceptNet. We extracted respective entries from a ConceptNet RDF dump[5] accessible via the MIT Media Lab ConceptNet homepage. We compiled a set of candidate queries by comparing the queries to the triples' 'RightText' segment. Queries and 'RightText' segments were pre-processed as described previously. If a given search query contained all tokens of the 'RightText' segment tokens, the query was added to the candidate set and the 'LeftText' segment was assigned the query as corresponding action. The following illustrative example should clarify the procedure:

|  | LeftText | Relation Type | RightText |
|---|---|---|---|
| ConceptNet Entry: | wait tables | MotivatedByGoal | make money |
| Potential New ConceptNet Entry: | wait tables | MotivatedByGoal | make some money quickly |

We conducted a human subject study to evaluate the compiled search query motivation/ ConceptNet action pairs. We had two independent annotators annotate a sampled set of 120 motivations. On the whole, 528 decisions had to be made. Table 11 displays an excerpt of the data to annotate. The annotators were given a list of motivations and corresponding actions. For every motivation/action pair (M/A), they had to answer following question with yes or no:

 *"Do you think that a person's motivation could be M when performing action A?"*

We introduced a softer variant of the known precision metric, i.e., a search query was considered a potential motivation if at least one ConceptNet action had been positively annotated. In our subject study, we achieved an average soft-precision of 64%, meaning that 77 out of 120 goals

---

from Search Query Logs were regarded reasonable motivations for the given actions. In addition, we calculated Cohen's kappa $\kappa$ – a metric reflecting the agreement amongst the annotators. The corresponding $\kappa$ – value in this study was 0.36 representing fair agreement [0.20 – 0.40]. Yet, we believe that our low $\kappa$-value partly results from the difficulty of assessing commonsense knowledge and the classification task itself. Both annotators had to decide whether an action/motivation pair was conceivable or reasonable. This subjective decision might be biased by each annotator's own experiences. Because of this bias, we applied a softer precision metric in our investigations.

**Table 5: Annotators had to decide whether the search queries in the left column could serve as motivation for every given action in the right column. Actions were automatically extracted from ConceptNet.**

| Goal from Search Query Logs | List of ConceptNet Actions |
|---|---|
| make some money quickly | wait tables, go to work, work the box office, serve customers, tell a story, get a contract, buy a house, apply for a job, pass a course |
| make new friends in your area | meet interesting people, meet people |
| find credit information | surf the net, surf the web, use a computer |
| ways to gain weight | eat ice cream |
| lose maximum weight fast | go jogging, eat healthly, release your energy, go for a run, play sports, get exercise, get some physical activity, eat vegetables |

Next, we review related literature relevant to our guiding research questions.

## RELATED WORK

In previous research, He et al. [6] have studied the acquisition of user goals from search *result snippets* (i.e. the segments of text listed on the result pages of search engines). Our work is novel in the sense that it studies search queries themselves as a source of common human goals, which can be suspected to better reflect intent. In the context of user intentions and query log analysis, search queries have been the subject of research for several years: An influential study conducted by Broder in 2000 [3] introduced a high level taxonomy of search intent, proposing a distinction between three classes of search goals: navigational, informational and transactional queries. This has stimulated a series of follow-up research on category refinement and automatic query categorization [2][8][10][21]. While previous research in this area has achieved considerable progress in the *categorization of queries* into high-level goal taxonomies serving a primarily *functional purpose* (to improve search, cf. [2][8][10]), we know little about the *acquisition of goal instances (common human goals)* from search query logs for *knowledge capture* purposes (as in [14][16]). Beyond related research in query log analysis

[3][8][10][17][18], our work is relevant to other research areas, notably commonsense knowledge acquisition and reasoning, goal mining and intelligent user interfaces.

Examples of related research in the area of *Commonsense Knowledge* are Openmind Commonsense and the related ConceptNet project, which tap into "knowledge contributions" made by volunteers on the web. ConceptNet contains a wide variety of commonsense knowledge including knowledge about human motivations and intentions [16]. ConceptNet's relation types "MotivationOf" and "DesireOf" are examples of attempts aiming to capture different aspects of intentional commonsense knowledge. In addition to volunteer-based knowledge acquisition, search engine result snippets have been studied for that purpose as well [5], and the idea of "human computation" [28] has inspired the development of games for collecting commonsense goals [14]. In relation to this body of work - to the best of our knowledge – our research represents the first *comparative* investigation of the utility of *search queries* for capturing common human goals.

Another area related to our research is *Goal Mining*, which is often referred to as the acquisition of goals from textual resources. This research area covers a range of interesting problems, including the acquisition of goals from patents, scientific articles [7], organizational policies [20], organizational guidelines and procedures [13] and others. To the best of our knowledge, acquiring common human goals *from search queries* has not been studied before in the context of goal mining.

## THREATS TO VALIDITY

While our study reveals some potentials of search query logs for knowledge acquisition, in this section we will discuss selected threats to validity [30] and limitations:

*Construct Validity*: The main construct we investigate in this research is the notion of *common human goals*. While we have used a definition of goals that is based on existing work and proved useful in previous similar contexts [26], some of the analyses conducted in this work might introduce a bias to our results. For example, due to the absence of explicit "goal" nodes in ConceptNet, we had to decide on a particular way of acquiring nodes that are likely to represent goals, and our approach depends, to some extent, on this decision. In addition, the goal matching algorithm adopted to assess overlap between goals in query logs and in ConceptNet neglects the sequence of tokens in goal statements. This might lead to incorrect matching pairs such as the hypothetical example [listen to the entertainer vs. entertain the listener]. However, by inspecting samples of matches, we could find few examples of this problem.

*Reliability*: The automatic classification approach for acquiring queries containing goals relies on existing tools such as the WEKA toolkit [29] and algorithms, which yields us to believe that reproducing our results is possible.

## CONCLUSIONS

Our work illustrated some potentials of Search Query Logs to help address the problems associated with capturing knowledge about common human goals on the web. In a departure from existing approaches focusing on volunteer-based, game-based or other forms of contribution, we present the results of applying an automatic classification approach to a renewable resource (search query logs) and results from a comparative study that introduce search query logs as a *feasible*, yet largely *untapped* resource for this task. Because query logs are a natural byproduct of human activity on the web, the costs associated with knowledge acquisition could be lower compared to other approaches requiring, for example, knowledge engineers. Quantitative and qualitative analyses of our results revealed that the goals acquired from search query logs in part represent commonsense knowledge contained in existing commonsense knowledge bases, and that they cover a vast range of topics and levels of granularity, which makes search query logs an interesting resource for addressing the goal coverage problem. Although the preliminary investigations presented are promising, more research is necessary. Ongoing work focuses on the development of more sophisticated approaches, the utilization of additional methods such as linguistic methods, and further evaluations to address the construction of conceptual nets representing and correlating broad sets of knowledge about common human goals.

## REFERENCES

[1] J. Allan and H. Raghavan. Using part-of-speech patterns to reduce query ambiguity. Proceedings of the 25th annual international ACM SIGIR Conference on Research and Development in Information Retrieval, 307--314, ACM Press New York, NY, USA,2002.

[2] R.A. Baeza-Yates, L. Calderon-Benavides and C.N. Gonzalez-Caro. The intention behind web queries. In F. Crestani and P. Ferragina and M. Sanderson, editor(s), Proceedings of String Processing and Information Retrieval (SPIRE), (4209):98--109, Springer, 2006.

[3] A. Broder, A taxonomy of web search, SIGIR Forum, vol. 36, no. 2, pp. 3-10, 2002.

[4] J. Cohen. A coefficient of agreement for nominal scales. Educational and Psychological Measurement, (20)1:37, 1960.

[5] I. S. Eslick. Searching for commonsense. Masters Thesis, MIT, 2006.

[6] K.Y. He, Y.S. Chang and W.H. Lu. Improving Identification of latent user goals through search-result snippet classification. In Proceedings of the 2007 IEEE/WIC/ACM International Conference on Web Intelligence, 683-686, IEEE Computer Society, 2007.

[7] B. Hui and E. Yu. Extracting conceptual relationships from specialized documents. Data & Knowledge Engineering, (54)1:29-55, Elsevier, 2005.

[8] B.J. Jansen, D.L. Booth and A. Spink. Determining the informational, navigational, and transactional intent of Web queries. Information Processing and Management, (44)3:1251--1266, Elsevier, 2008.

[9] D. Kirsh. When is information explicitly represented? Information, Language and Cognition - The Vancouver Studies in Cognitive Science., 340-365, UBC Press,1990.

[10] U. Lee, Z. Liu and J. Cho. Automatic identification of user goals in Web search. In Proceedings of the 14th International World Wide Web Conference (WWW'05), 391-400, ACM Press,New York, NY, USA,2005.

[11] D.B. Lenat. CYC: A large-scale investment in knowledge infrastructure. Communications of the ACM, (38)11:33-38, 1995.

[12] B. Levin, English verb classes and alternations: a preliminary investigation, University of Chicago Press, 1993.

[13] S. Liaskos, A. Lapouchnian, Y. Yu, E. Yu and J. Mylopoulos. On goal-based variability acquisition and analysis. In Proceedings of the 14th IEEE International Requirements Engineering Conference (RE'06), Minneapolis, USA, 2006.

[14] H. Lieberman, D.A. Smith and A. Teeters. Common Consensus: a web-based game for collecting commonsense goals. In Proceedings of the Workshop on Common Sense and Intelligent User Interfaces held in conjunction with the 2007 International Conference on Intelligent User Interfaces (IUI 2007), 2007.

[15] H. Liu, H. Lieberman and T. Selker. GOOSE: A goal-oriented search engine with commonsense. AH '02: Proceedings of the Second International Conference on Adaptive Hypermedia and Adaptive Web-Based Systems, 253--263, Springer-Verlag,London, UK,2002.

[16] H. Liu and P. Singh. ConceptNet - A practical commonsense reasoning tool-kit. BT Technology Journal, (22)4:211-226, 2004.

[17] G.C. Murray and J. Teevan. WWW Workshop Report: Query log analysis - social and technological challenges. ACM SIG IR Forum, (41)2 ACM New York, NY, USA, 2007.

[18] M. Pasca, B. Van Durme and N. Garera. The role of documents vs. queries in extracting class attributes from text. In Proceedings of the International Conference on Information and Knowledge Management (CIKM'07), 485-494, ACM, New York, NY, USA, 2007.

[19] G. Pass, A. Chowdhury and C. Torgeson. A picture of search. Proceedings of the 1st International Conference on Scalable Information Systems, ACM Press New York, NY, USA, 2006.

[20] C. Potts, K. Takahashi and A. I. Anton. Inquiry-based requirements analysis. IEEE Software, (11)2, 1994.

[21] D. E. Rose and D. Levinson, Understanding user goals in web search. In Proc. of WWW 2004, May 17-22, 2004, New York, USA, 2004.

[22] G. Regev and A. Wegmann. Where do goals Come from: the underlying principles of goal-oriented requirements engineering. RE '05: Proceedings of the 13th IEEE International Conference on Requirements Engineering (RE'05), 253--362, IEEE Computer Society, Washington, DC, USA, 2005.

[23] R.C. Schank and R.P. Abelson. Scripts, plans, goals, and understanding: an inquiry into human knowledge structures. Lawrence Erlbaum Associates, 1977.

[24] P. Singh, T. Lin, E.T. Mueller, G. Lim, T. Perkins and W.L. Zhu. Open Mind Common Sense: Knowledge acquisition from the general public. In Proceedings of the First International Conference on Ontologies, Databases, and Applications of Semantics for Large Scale Information Systems, 1223--1237, Springer-Verlag London, UK, 2002.

[25] D.A. Smith. EventMinder: A Personal Calendar Assistant That Understands Events. Master Thesis, MIT, 2007.

[26] M. Strohmaier, P. Prettenhofer and M. Kröll: Acquiring Explicit User Goals From Search Query Logs, *in* 'Internat. Workshop on Agents and Data Mining Interaction ADMI ' 08, in conjunction with WI '08', 2008.

[27] M. Strohmaier, P. Prettenhofer and M. Lux. Different degrees of explicitness in intentional artifacts - studying user goals in a large search query log. Proceedings of the CSKGOI'08 Workshop on Commonsense Knowledge and Goal Oriented Interfaces, held in conjunction with IUI'08, Canary Islands, Spain, 2008.

[28] L. von Ahn. Games with a purpose. Computer, (39)6:92--94, IEEE Computer Society Press Los Alamitos, CA, USA, 2006.

[29] I. H. Witten and E. Frank, Data Mining: Practical machine learning tools and techniques, 2nd ed., ser. Morgan Kaufmann Series in Data Management Systems. Morgan Kaufmann, June 2005.

[30] R. K. Yin, Case study research: design and methods (Applied Social Research Methods). SAGE Publications, December 2002.