

# Overview of the NTCIR-11 IMine Task

Yiqun Liu<sup>1</sup>, Ruihua Song<sup>2</sup>, Min Zhang<sup>1</sup>, Zhicheng Dou<sup>2</sup>, Takehiro Yamamoto<sup>3</sup>,  
Makoto Kato<sup>3</sup>, Hiroaki Ohshima<sup>3</sup>, Ke Zhou<sup>4</sup>

<sup>1</sup>Tsinghua University, <sup>2</sup>Microsoft Research Asia, <sup>3</sup>Kyoto University, <sup>4</sup>University of Edinburgh  
{yiqunliu, z-m}@tsinghua.edu.cn; {song.ruihua, zhichdou@microsoft.com}; {tyamamot, kato,  
ohshima@dl.kuis.kyoto-u.ac.jp}; zhouke.nlp@gmail.com

## ABSTRACT

In this paper, we provide an overview of the NTCIR IMine task, which is a core task of NTCIR-11 and also a succeeding work of INTENT@NTCIR-9 and INTENT2@NTCIR-10 tasks. IMine is composed of a subtopic mining (SM) task, a document ranking (DR) task and a TaskMine (TM) pilot task. 21 groups from Canada, China, Germany, France, Japan, Korea, Spain, UK and United States registered to the task, which makes it one of the largest tasks in NTCIR-11. Finally, we receive 45 runs from 10 teams to the SM task and 25 runs from 6 groups to the DR task. We describe the task details, annotation of results, evaluation strategies and then the official evaluation results for each subtask.

## Keywords

Intent, ambiguity, diversity, evaluation, test collection.

## 1. INTRODUCTION

Many queries are short and vague in practical Web search environment. By submitting one query, users may have different intents. For an ambiguous query, users may seek for different interpretations. For a query on a broad topic, users may be interested in different subtopics. Today mining users' underlying intents of a query is an interesting topic for both IR communities and commercial search engines. IMine task aims to provide common data sets and evaluation methodology to researchers who want to investigate into the techniques for better understanding user intents behind ambiguous or broad queries. IMine is short for search Intent Mining and it also pronounces like “暧昧” which means “ambiguous” in Chinese and Japanese.

Through IMine task, we expect participants to advance the state-of-the-art techniques explored in INTENT [1] and INTENT2 [2] and to gain further insight into the right balance between relevance and diversity. We involve more user behavior data both for participants and in the annotation process to help assessors for subtopic clustering and importance estimation. We are also interested in comparing the differences between diversified search annotations from a small number of professional assessors and a relatively large number of untrained users as crowd sourcing efforts.

Similar with INTENT tasks, the IMine task consists of two subtasks: Subtopic Mining and Document Ranking. While the SM task may be regarded as a pre-DR task for identifying explicit intents, it can also be useful for other practical tasks such as query suggestion and auto-completion. We also setup a pilot subtask named TaskMine which focus on exploiting the techniques of understanding the relationship among tasks for supporting the Web searchers. We involve dealing with three different languages including English, Chinese and Japanese in IMine task. Query topics for all three languages were developed for SM task while only English and Chinese DR tasks are required since few participants show interests in Japanese DR. The major differences between IMine and previous INTENT2 tasks are shown in Table 1.

Table 1. Differences between IMine and INTENT2 tasks

|                              | INTENT2  | IMINE   |
|------------------------------|--|---|
| <b>Number of Topics</b>      | Chinese: 100<br>Japanese: 100<br>English: 50         | Chinese: 50<br>Japanese: 50<br>English: 50  |
| <b>DR task corpus</b>        | Chinese: SogouT<br>Japanese:<br>ClueWeb JA           | Chinese: SogouT<br>English:<br>ClueWeb12-B13  |
| <b>Crowd sourcing</b>        | No   | Crowd sourcing or Chinese DR  |
| <b>Subtopic organization</b> | One level  | Two level: no more than 5 first-level subtopics with at most 10 second-level subtopics each   |
| <b>Subtopic candidate</b>    | Query suggestions from Bing, Google, Sogou and Baidu | Query suggestions from Bing, Google, Sogou, Yahoo! and Baidu; Query facets generated by [3] from search engine results; Query facets generated by [4] from Sogou log data |
| <b>User behavior data</b>    | SogouQ (data collected in 2008): appr. 2GB           | SogouQ (data collected in 2008 and 2011): appr. 4GB   |

From Table 1 we can see that there are two major differences between IMine and previous INTENT tasks. The first difference lies that IMine requires participants to submit a two-level hierarchy of sub-intents for the query topics. In previous diversified search related studies, we notice the phenomena that some query subtopics belong to the concept of others (e.g. *iPhone* and *apple inc. products* are both regarded as subtopics for the query *apple*, while *iPhone* should be covered by *apple inc. products*). This may lead to difficulty in subtopic importance estimation and diversified ranking. Therefore, we introduce a two-level hierarchy of subtopics to better present the diversified intent structure of ambiguous/broad queries. This require extra efforts in assessment and a different design of evaluation metrics, which we will address in follow up sections.

The second major difference between IMine and previous tasks is that we try to incorporate more user behavior data and introduce the evaluation framework based on crowd sourcing. Recently, several metrics have been proposed to evaluate a diversified search result with different types of user behavior assumptions, considering relevance, diversity, novelty, user intent, and so on. To validate the credibility of these evaluation metrics, a number of methods that "evaluate evaluation metrics" are also adopted in diversified search evaluation studies, such as Kendall's tau [5], Discriminative Power [6], and the Intuitiveness Test [7]. These methods have been widely

adopted and have aided us in gaining much insight into the effectiveness of evaluation metrics. However, they also follow certain types of user behaviors or statistical assumptions and do not take the information of users' actual search preferences into consideration. In IMine task, we want to take user preferences collected with crowd sourcing efforts as the ground truth to investigate into both the performance of participants' runs and diversified evaluation metrics.

21 groups from Canada, China, Germany, France, Japan, Korea, Spain, UK and United States registered to the IMine task, which makes it one of the largest tasks in NTCIR-11. Finally, we receive 45 runs from 10 groups to the SM task and 25 runs from 6 groups to the DR task. Names and organizations of the participants which submitted results are shown in Table 3 and Table 4.

**Table 3. Organization of the participating groups in IMine**

| Group Name | Organization  |
|------------|---|
| UDEL       | University of Delaware, United States                                   |
| SEM13      | Toyohashi University of Technology, Japan                               |
| HULTECH    | University of Caen, France  |
| THU-SAM    | Joint team of Tsinghua University, China and Samsung Electronics, Korea |
| FRDC       | Fujitsu Research & Development Center Co., LTD., China                  |
| TUTA1      | The University of Tokushima, Japan                                      |
| CNU        | Capital Normal University, China  |
| KUIDL      | Kyoto University, Japan   |
| UM13       | University of Montreal, Canada  |
| KLE        | POSTECH, Korea  |

**Table 4. Result submission from different participating groups in IMine**

| Group         | Chinese SM | Japanese SM | English SM | Chinese DR | English DR |
|---------------|------------|-------------|------------|------------|------------|
| UDEL          |            |             | 1          |            | 5          |
| SEM13         |            |             | 5          |            | 5          |
| HULTECH       |            |             | 4          |            |            |
| THU-SAM       | 5          |             | 2          | 4          |            |
| FRDC          | 5          |             |            | 5          |            |
| TUTA1         | 1          |             | 1          | 1          | 2          |
| CNU           | 4          |             |            |            |            |
| KUIDL         |            | 1           | 1          |            |            |
| UM13          |            |             | 3          |            | 3          |
| KLE           | 4          | 4           | 4          |            |            |
| <b>#Group</b> | <b>5</b>   | <b>2</b>    | <b>8</b>   | <b>3</b>   | <b>4</b>   |
| <b>#Run</b>   | <b>19</b>  | <b>5</b>    | <b>21</b>  | <b>10</b>  | <b>15</b>  |

The remainder of the paper is organized as follows: Section 2 describes the details of the three subtasks, including the query set, supporting data resources and the test corpus adopted. The evaluation metrics and result assessment process are introduced in Section 3. Official evaluation results based on cranfield methodology are presented in Section 4. User preference test results are reported and compared with cranfield-like approaches in Section 5. Section 6 concludes this paper and the Appendix contains the details of each run as well as significance test results.

## 2. TASKS AND DATASETS

### 2.1 Query set

The same query topics are adopted in both Subtopic Mining and Document Ranking subtasks for all languages. These topics are sampled from the median-frequency queries collected from both Sogou and Bing search logs. We avoid top or tail queries because

search performance of top queries are already quite high for most commercial search engines while many tail queries may contain typos, language mistakes or even illegal contents. Approximately equal amounts of ambiguous, broad and clear queries are included in the query topic set. Several topics are shared among different languages for possible future cross-language research purposes. Detailed information of the constructed query set is shown in Table 5. For SM task, queries with clear intents are not evaluated because they are not expected to contain subtopics.

**Table 5. Statistics of the IMine query topic set**

| Language | #topic    |       |       | #shared topics  |
|----------|-----------|-------|-------|---|
|          | Ambiguous | Broad | Clear |   |
| English  | 16        | 17    | 17    | 14 shared by all languages, 8 shared by English and Chinese |
| Chinese  | 16        | 17    | 17    |   |
| Japanese | 17        | 17    | 16    |   |

We follow the query intent classification framework proposed in [8] and group the queries into three groups: Ambiguous, Broad and Clear. Both ambiguous and broad queries are adopted in the SM task for query intent analysis while all queries are evaluated in the DR task (for clear queries, we just evaluate the ad-hoc retrieval performance instead of diversified search performance).

**Table 6. IMine query topic set (for Intent, a: ambiguous, b: broad, c: clear)**

| ID   | Topic | Intent | Shared |
|------|-------|--------|--------|
| 0001 | 先知    | a      | CEJ    |
| 0002 | 波斯猫   | a      | CE     |
| 0003 | 猫头鹰   | a      | CEJ    |
| 0004 | Adobe | a      | CEJ    |
| 0005 | 传奇    | a      | CEJ    |
| 0006 | 小米    | a      |        |
| 0007 | 中国水电  | a      |        |
| 0008 | 云轩    | a      |        |
| 0009 | 遮天    | a      |        |
| 0010 | 舍得    | a      |        |
| 0011 | 秋菊    | a      |        |
| 0012 | 三字经   | a      |        |
| 0013 | 三毛    | a      |        |
| 0014 | 阳光    | a      |        |
| 0015 | 嫦娥    | a      |        |
| 0016 | 程序员   | a      |        |
| 0017 | 泰国特产  | b      | CE     |
| 0018 | 科学美国人 | b      | CEJ    |
| 0019 | 黄金    | b      | CE     |
| 0020 | 浴缸    | b      | CEJ    |
| 0021 | 婚戒    | b      | CEJ    |
| 0022 | 三星    | b      | CEJ    |
| 0023 | 饥饿游戏  | b      | CEJ    |
| 0024 | 心理测试  | b      |        |
| 0025 | 椰岛造型  | b      |        |
| 0026 | 野葛根   | b      |        |
| 0027 | 秧歌    | b      |        |
| 0028 | 卫子夫   | b      |        |
| 0029 | 佛教音乐  | b      |        |
| 0030 | 浏览器下载 | b      |        |

|      |  |   |     |
|------|--|---|-----|
| 0031 | 相亲节目有哪些  | b |     |
| 0032 | 哈利波特   | b |     |
| 0033 | 安卓 2.3 游戏下载                                    | b |     |
| 0034 | 男鞋尺码对照表  | c | CEJ |
| 0035 | 奥巴马简历  | c | CE  |
| 0036 | 肥胖的原因  | c | CEJ |
| 0037 | 什么是自然数   | c | CEJ |
| 0038 | 牙齿黄怎么办   | c | CE  |
| 0039 | 治疗近视的方法  | c | CE  |
| 0040 | 央金兰泽的歌曲  | c |     |
| 0041 | 声卡是什么  | c |     |
| 0042 | 乘法口诀   | c |     |
| 0043 | 学雷锋作文  | c |     |
| 0044 | 联通网上营业厅  | c |     |
| 0045 | 怎么查 ip 地址                                      | c | CEJ |
| 0046 | 邮编号码查询   | c | CEJ |
| 0047 | 在线冲印照片   | c | CE  |
| 0048 | qq 加速器下载                                       | c |     |
| 0049 | 冬季恋歌国语全集                                       | c |     |
| 0050 | 初恋这件小事   | c |     |
| 0051 | apple  | a |     |
| 0052 | cathedral                                      | a |     |
| 0053 | eclipse  | a |     |
| 0054 | fas  | a |     |
| 0055 | flesh  | a |     |
| 0056 | ir   | a |     |
| 0057 | lost   | a |     |
| 0058 | shrew  | a |     |
| 0059 | symmetry                                       | a |     |
| 0060 | the presidents of the united states of america | a |     |
| 0061 | windows  | a |     |
| 0062 | prophet  | a | CEJ |
| 0063 | gold   | a | CE  |
| 0064 | owl  | a | CEJ |
| 0065 | adobe  | a | CEJ |
| 0066 | legend   | a | CEJ |
| 0067 | beijing subways                                | b |     |
| 0068 | camera   | b |     |
| 0069 | free dvd burner                                | b |     |
| 0070 | lost season 5                                  | b |     |
| 0071 | mobile phones                                  | b |     |
| 0072 | programming languages                          | b |     |
| 0073 | tom cruise                                     | b |     |
| 0074 | top ipad games                                 | b |     |
| 0075 | watches  | b |     |
| 0076 | thai specialties                               | b | CE  |
| 0077 | scientific american                            | b | CEJ |
| 0078 | persian cat                                    | b | CE  |
| 0079 | bathtub  | b | CEJ |
| 0080 | wedding ring                                   | b | CEJ |
| 0081 | samsung  | b | CEJ |
| 0082 | the hunger games                               | b | CEJ |
| 0083 | harry potter                                   | b | CE  |
| 0084 | 21 weeks pregnant                              | c |     |
| 0085 | 7zip   | c |     |
| 0086 | appendix pain symptoms                         | c |     |
| 0087 | brad paisley lyrics                            | c |     |

|      |                             |   |     |
|------|-----------------------------|---|-----|
| 0088 | craig's list phoenix        | c |     |
| 0089 | mcdonalds nutrition guide   | c |     |
| 0090 | sausalito art festival      | c |     |
| 0091 | tennessee unemployment      | c |     |
| 0092 | men's shoe sizes conversion | c | CEJ |
| 0093 | obama biography             | c | CE  |
| 0094 | causes of obesity           | c | CEJ |
| 0095 | what is a natural number    | c | CEJ |
| 0096 | yellow teeth treatment      | c | CE  |
| 0097 | myopia treatment            | c | CE  |
| 0098 | how to find my ip address   | c | CEJ |
| 0099 | postcode finder             | c | CEJ |
| 0100 | online photo printing       | c | CE  |
| 0101 | シド                          | a |     |
| 0102 | ダム                          | a |     |
| 0103 | R                           | a |     |
| 0104 | ハヤブサ                        | a |     |
| 0105 | ナポレオン                       | a |     |
| 0106 | アバター                        | a |     |
| 0107 | ジップ                         | a |     |
| 0108 | ウォッカ                        | a |     |
| 0109 | 横浜                          | a |     |
| 0110 | 伝奇                          | a | CEJ |
| 0111 | アドビ                         | a | CEJ |
| 0112 | 予言者                         | a | CEJ |
| 0113 | オウル                         | a | CEJ |
| 0114 | 赤とうがらし                      | a |     |
| 0115 | 銀シャリ                        | a |     |
| 0116 | 嵐                           | a |     |
| 0117 | フランクフルト                     | a |     |
| 0118 | 東方神起                        | b |     |
| 0119 | 円形脱毛症                       | b |     |
| 0120 | 柿の葉すし                       | b |     |
| 0121 | シャネル                        | b |     |
| 0122 | 女子バレー                       | b |     |
| 0123 | TPP                         | b |     |
| 0124 | ドラえもん                       | b |     |
| 0125 | ビートルズ                       | b |     |
| 0126 | ボーカロイド                      | b |     |
| 0127 | 年賀状                         | b |     |
| 0128 | うつ病                         | b |     |
| 0129 | サムスン                        | b | CEJ |
| 0130 | タイ 特産                       | b | CEJ |
| 0131 | 浴槽                          | b | CEJ |
| 0132 | ハンガーゲーム                     | b | CEJ |
| 0133 | 結婚指輪                        | b | CEJ |
| 0134 | サイエンティフィック・アメリカン            | b | CEJ |
| 0135 | 櫻井歯科診療所 ホームページ              | c |     |
| 0136 | 京葉タクシー 電話番号                 | c |     |
| 0137 | 湘南新宿ライン 路線図                 | c |     |
| 0138 | 旭山動物園 アクセス                  | c |     |
| 0139 | 秋田中央交通 時刻表                  | c |     |
| 0140 | 羽田空港 リムジンバス 時刻表             | c |     |

|      |                   |   |     |
|------|-------------------|---|-----|
| 0141 | 水平投射運動 速度の求め方     | c |     |
| 0142 | のし袋 書き方           | c |     |
| 0143 | facebook 退会方法     | c |     |
| 0144 | タロットカード 吊るされた男 意味 | c |     |
| 0145 | 少々お待ちください 英語      | c |     |
| 0146 | 肥満の原因             | c | CEJ |
| 0147 | 自然数とは             | c | CEJ |
| 0148 | IP アドレスを確認するには    | c | CEJ |
| 0149 | メンズ靴サイズ対応表        | c | CEJ |
| 0150 | 郵便番号検索            | c | CEJ |

## 2.2 Subtopic Mining Subtask

In the Subtopic Mining task, a subtopic could be an interpretation of an ambiguous query or an aspect of a broad query. Participants are expected to generate a two-level hierarchy of underlying subtopics by analysis into the provided document collection, user behavior data set or other kinds of external data sources. A list of query suggestions/completions collected from popular commercial search engines as well as some queries mined from search logs/SERPs (see Table 1) are provided as possible subtopic candidates while participants can also use other information sources (e.g. Wikipedia, search behavior logs) to generate their own candidates. For both Subtopic Mining and Document Ranking subtasks, SogouQ search user behavior data collection is available for participants as additional resources. The collection contains queries and click-through data collected and sampled by China's second largest search engine Sogou.com in 2008 and 2012, separately.

As for the two-level hierarchy of subtopics, we can take the ambiguous query “windows” as an example. The first-level subtopic may be Microsoft Windows, software in windows platform or house windows. In the category of Microsoft Windows, users may be interested in different aspects (second-level subtopics), such as “Windows 8”, “Windows update”, etc.

From our experiences in past INTENT/INTENT2 tasks, we found that the relationship among subtopics for some queries are not trivial. For example, for topic #0205 in INTENT2 task (功夫/kung fu), the subtopics 功夫【电影《功夫》】(kung fu the movie), 功夫【影片下载】(movie download), 功夫【在线观看】(movie online), 功夫【影视作品】(other movies related with kungfu) should be grouped into a same category “movies related with kungfu” instead of several different subtopics. We believe that organizing such a hierarchical structure of subtopics will help search engines to present a better ranking of results.

The hierarchical structure of subtopics is close related with knowledge graph which has been well studied in Web search researches recently. Some participants in INTENT/INTENT2 tasks also adopted existing knowledge graphs such as wikipedia, freebase (e.g. THCI and THUIS in INTENT2) in developing subtopic candidate sets. However, we believe that the hierarchical subtopics for a certain query is used to describe users' possible information needs behind this query instead of the knowledge structure of the entity named this query. Therefore, even when a knowledge graph exists for a given query (which is not usually the case since Web queries are so complicated), we should not use the

graph directly as the hierarchy of query intents.

In this year's IMine task, at most FIVE first-level subtopics with no more than TEN second-level subtopics each should be returned for each query topic. There is no need to return subtopics for clear queries but participants will not be penalized for doing this in the evaluation. The first-level subtopics for broad queries will not be taken into consideration in the evaluation process because there may be various standards for organizing high-level aspects for these queries. Besides the hierarchy of subtopics, a ranking list of all first-level subtopics and a separate ranking list of all second-level subtopics should also be returned for each ambiguous/broad query. It means that the submitted second-level subtopics should be globally ranked across different first-level subtopics within a same query topic. With these ranking lists, the importance estimation results could be evaluated and compared among different participant runs.

## 2.3 Document Ranking Subtask

In document ranking task, Participants are asked to return a diversified ranked list of no more than 100 results for each query. Participants are encouraged to selectively use diversification algorithms in ranking because diversification is not necessary for all queries (e.g. for clear queries). Based on the subtopic mining results, participants are supposed to select important first-level/second-level subtopics and mix them to form a diversified ranking list. The goals of diversification are (a) to retrieve documents that cover as many intents as possible; and (b) to rank documents that are highly relevant to more popular intents higher than those that are marginally relevant to less popular intents.

SogouT (<http://www.sogou.com/labs/dl/t-e.html>) is adopted as the document collection for Chinese topics in Document Ranking subtask. The collection contains about 130M Chinese pages together with the corresponding link graph. The size is roughly 5TB uncompressed. The data was crawled and released on Nov 2008. In order to help participants who are not able to construct their own retrieval platforms, the organizers provide a non-diversified baseline Chinese DR run based on THU-SAM's retrieval system.

As for English Document Ranking subtask, the ClueWeb12-B13 (<http://lemurproject.org/clueweb12/>) data set is adopted, which includes 52M English Web pages crawled in 2012. A search interface is provided by Lemur project so that the retrieval baseline could be obtained without having to construct one's own search index.

## 2.4 TaskMine Subtask

//to be added by TaskMine organizers

# 3. EVALUATION METRICS

## 3.1 Subtopic Mining and Document Ranking Subtasks

Search result evaluations are based on the document relevance assessments with respect to certain queries. Supposing that these documents are assessed with level 0 to  $h$  where 0 means irrelevant and  $h$  means the highest relevant. Hence  $h=1$  means a binary relevance assessment. Let  $N_x$  denote the number of relevant documents at level  $x$  ( $0 < x < h$ ), then  $N = \sum_x N_x$  means the total

number of relevant documents. Let  $d_r$  denote the document at rank  $r$  in the result list and define  $J(r)=1$  if  $d_r$  is relevant to a query at level  $x$  ( $0 < x < h$ ), otherwise  $J(r)=0$ . We denote the cumulative

number of relevant documents as  $C(r) = \sum_{i=1}^r J(i)$ .

Let  $g(r)$  denote the document gain of  $d_r$ , then  $cg(r) = \sum_{i=1}^r g(i)$  means the cumulative gain at rank  $r$ . Also, the gain and cumulative gain of the ideal ranked list are denoted as  $g^*(r)$  and  $cg^*(r)$  respectively. Then we can define  $nDCG$  at document cutoff  $l$  as:

$$nDCG@l = \frac{\sum_{r=1}^l g(r) / \log(r+1)}{\sum_{r=1}^l g^*(r) / \log(r+1)} \quad (1)$$

Diversified search evaluation requires document relevance assessments with respect to subtopics instead of queries, which is different from the traditional evaluations. Document gains are therefore evaluated in terms of subtopics underlying the query. Let  $g_i(r)$  denote the gain of  $d_r$  with respect to subtopic  $i$ ,  $N_i$  denote the total number of documents relevant to subtopic  $i$ , and  $J_i(r)$  indicate whether  $d_r$  is relevant to subtopic  $i$ . Furthermore, we suppose that there are  $n$  subtopics underlying a query  $q$  and denote the probability distribution of subtopic  $i$  as  $P(i|q)$ , therefore  $\sum_{i=1}^n P(i|q) = 1$ .

In INTENT/INTENT2 tasks, the major evaluation metric is  $D\#-measures$  which is proposed in [9] to more intuitively evaluate the diversity of a ranked list. The main idea is that the abandonment of the separate calculation of measures for each subtopic, which is leveraged in previous IA measures proposed in [10] and [11]. By introducing a new document gain (named *Global Gain*), the original document gains calculated in terms of each subtopic are linearly combined. The *Global Gain* is defined as follows:

$$GG(r) = \sum_{i=1}^n P(i|q) g_i(r) \quad (2)$$

Then document gains in the traditional measures are replaced by this *Global Gain* factor. After this replacement, these measures (referred to as  $D\#-measures$ ) capture all the properties of the original measures. Furthermore, the *Global Gain* linearly combines the original document gain with the respective subtopic probability for each document in an overall perspective, which directly reflects the diversity. To evaluate the subtopic recall, [9] also defined the measure namely  $I-rec$ <sup>1</sup>, which is the proportion of subtopics covered by documents:

$$I-rec@l = |\bigcup_{r=1}^l I(r)| / n \quad (3)$$

where  $I(r)$  stands for the set of subtopics which  $d_r$  is relevant to. Linearly combining the  $D\#-measures$  with  $I-rec$  for documents at cutoff  $l$ , [9] defined the  $D\#-measures$  as follows:

$$D\#-measure@l = \lambda I-rec@l + (1-\lambda) D\#-measure@l \quad (4)$$

where  $\lambda$  is the tradeoff between the diversity and the subtopic recall and is set to 0.5 in [12]. The  $D\#-measures$  are adopted in both subtopic mining and document ranking tasks in INTENT/INTENT2 tasks as the main evaluation metric.

Besides the  $D\#-measures$ ,  $DIN-measures$  were also adopted in INTENT2 task. According to [14], diversity evaluation should distinguish the navigational subtopic from the informational one. The reason lies that when a certain subtopic is a navigational one, the user wants to see only one particular web page; while the user is happy to see many relevant pages when the subtopic is informational. Therefore, the types of information needs behind subtopics should be taken into account and different measures should be leveraged for evaluating subtopics in different types.

Based on this assumption, the reformulation of the *Global Gain* factor in  $DIN-measures$  is described as follows:

$$GG^{DIN}(r) = \sum_i P(i|q) g_i(r) + \sum_j isnew_j(r) P(j|q) g_j(r) \quad (5)$$

where  $\{i\}$  and  $\{j\}$  denote the sets of informational and navigational subtopics for query  $q$ . And  $isnew_j(r)$  is an indicator that if there is no document relevant to the navigational subtopic  $j$  between ranks 1 and  $r-1$ ,  $isnew_j(r)$  is set to 1, otherwise  $isnew_j(r)$  is set to 0. In this way,  $GG^{DIN}$  evaluates the informational and navigational subtopics in different ways. From this definition, we can find that  $GG^{DIN}$  evaluate the informational subtopic in the same way as  $D\#-measures$ , but for the navigational subtopic  $j$ , it leverages the indicator  $isnew_j(r)$  to guarantee that only the first relevant document is considered. The  $DIN-measures$  are then calculated by replacing the  $GG(r)$  of  $D\#-measures$  with  $GG^{DIN}$ .

In IMine task, we follow the settings in INTENT/INTENT2 and choose  $D\#-nDCG$  as the main evaluation metric for Document Ranking subtask. However, since a hierarchy instead of a single list of subtopics are submitted for each query topic in the new Subtopic Mining task, new metrics should be designed to evaluate the performance of the submitted two-level hierarchy of subtopics.

For the IMine Subtopic Mining task, we propose to use the  $H-measures$  (evaluation measures of Hierarchical subtopic structure) as the main evaluation metric. The definition of  $H$ -measure is as follows:

$$H - measure = Hscore * (\alpha * Fscore + \beta * Sscore), \quad (\alpha + \beta = 1) \quad (6)$$

The definitions of  $Hscore$ ,  $Fscore$  and  $Sscore$  are as follows and they each describe one aspect of the submitted hierarchy.

$Hscore$  measures the quality of the hierarchical structure by whether the second-level subtopic is correctly assigned to the appropriate first-level subtopic.

$$Hscore = \frac{\sum_{i=1}^{N^{(1)}} accuracy(i)}{N^{(1)}} \quad (7)$$

Here  $N^{(1)}$  is the number of first-level subtopics for a certain query topic in the submission (no more than 5).  $accuracy(i)$  is the percentage of correctly-assigned second-level subtopics for first-level subtopic  $i$ . If first-level subtopic  $i$  is not relevant to the query topic, then  $accuracy(i)$  should be 0. Irrelevant second-level subtopics should not be regarded as “correctly-assigned” ones.

$Fscore$  measures the quality of the first-level subtopic by whether the submitted first-level subtopics are correctly ranked and whether all important first-level subtopics are found:

$$Fscore = D\#-measure(FS_1, FS_2, \dots, FS_{N^{(1)}}) \quad (8)$$

Here  $\{FS_i\}$  is the first-level subtopic list for a certain query topic ranked by the score contained in submission file.

Similar with  $Fscore$ ,  $Sscore$  measures the quality of the second-level subtopic with the following equation:

$$Sscore = D\#-measure(SS_1, SS_2, \dots, SS_{N^{(2)}}) \quad (9)$$

Here  $\{SS_i\}$  is the second-level subtopic list for a certain query topic ranked by multiplying the scores of the second-level subtopic and its corresponding first-level subtopic. Notice that all second-level subtopics are globally ranked in the submitted results so that a single  $\{SS_i\}$  list could be derived.

<sup>1</sup> In [12] the authors renamed the  $S-Recall$  proposed in [13] as  $I-rec$ .

We can see that the parameters  $\alpha$  and  $\beta$  are used to balance the scores of first-level and second-level subtopics. Note that the first level subtopics are not considered in the evaluation of broad queries because there may be different categories to group the second level subtopics (e.g. book/character/film or secret chamber/order of phoenix/death hollow for the query harry potter). Therefore,  $\alpha$  is set to 0 for all broad queries. As for ambiguous queries, we choose equal values of  $\alpha$  and  $\beta$  ( $\alpha = \beta = 0.5$ ).

### 3.2 TaskMine Subtask

//to be added by TaskMine organizers

## 4. RESULT ASSESSMENT

The result assessment process is completed by different groups of assessors. As for the Chinese and English SM/DR subtasks, a vendor company is hired by NII to finish the annotation. Meanwhile, assessment of the Japanese SM task is completed by volunteers recruited in Kyoto University. All the annotation tasks are completed by native speakers to guarantee quality.

### 4.1 Subtopic Mining Subtask

For the subtopic mining subtask, each ambiguous and broad query should be annotated by assessors to get a two-level hierarchy of subtopics. Clear queries are not considered in this subtask. The annotation process is completed in the following steps:

- Result pool construction: Result pool of the Chinese SM task contains 1,630 first-level subtopics, 6,594 second-level subtopics and 13,251 subtopic pairs (each pair is composed of a first-level subtopic and a corresponding second-level one as submitted by participating groups). Result pool of the Japanese SM task contains 539 first-level subtopics, 3,500 second-level subtopics and 5,467 subtopic pairs. Result pool of the English SM task contains 2,537 first-level subtopics, 13,993 second-level subtopics and 23,981 subtopic pairs.
- Annotation task 1 (relevance judgment): for each submitted first-level and second-level subtopic, the assessors are required to decide whether it is relevant to the query topic or not. Any irrelevant ones will be removed from the result pool and not dealt with in the following annotation tasks.
- Annotation task 2 (Subtopic relationship verification): For each second-level subtopic in a submission, the assessors are required to decide whether the submission correctly assigns its first-level subtopic.
- Annotation task 3 (first-level clustering): For all submitted hierarchy of subtopics, the assessors are required to cluster all the first-level subtopics into several clusters.
- Annotation task 4 (importance voting for first-level): For all first-level clusters, the assessors are required to vote for its importance and select the FIVE most important ones.
- Annotation task 5 (post-clustering classification): For all second-level subtopics, the assessors are required to decide which of the five most important first-level subtopic cluster it should belong to or it doesn't fit for any. The second-level subtopics that are not relevant to any first-level subtopic should be regarded as irrelevant.
- Annotation task 6 (second-level clustering): For each of the five most important first-level subtopics, the assessors are required to cluster all the second-level subtopics which belong to it into several clusters.
- Annotation task 7 (importance voting for second-level): For

all second-level clusters, the assessors are required to vote for its importance and retain at most TEN ones for each first-level cluster. The importance voting is for all second-level subtopics of the corresponding query instead of particular first-level subtopics.

With the above procedure, the two-level hierarchy of subtopics could be generated for each ambiguous/broad query topics. *Hscore* could be estimated with the results from Annotation task 2. Meanwhile *Fscore* and *Sscore* are estimated with the results generated in Annotation task 4 and Annotation task 7, separately. Note that in the calculation of *Hscore*, we do not consider whether the first-level or second-level subtopics are finally chosen as qrels or not. Instead, we want to evaluate whether the submitted hierarchy is self-consistent.

According to the assessment results for SM task, we have 116 first-level subtopics and 501 second-level subtopics for the 33 unclear queries in Chinese SM (3.51 first-level subtopics per query and 4.32 second-level subtopics per first-level subtopic on average). In English SM, we have 125 first-level subtopics and 373 second-level subtopics for the 33 unclear queries (3.79 first-level subtopics per query and 2.98 second-level subtopics per first-level subtopic on average). In Japanese SM, we have 145 first-level subtopics and 477 second-level subtopics for the 34 unclear queries (4.26 first-level subtopics per query and 3.29 second-level subtopics per first-level subtopic on average).

Although the participants are required to submit up to 10 second-level subtopics for each first-level subtopic, the assessment shows a much smaller number of second-level subtopics. We believe that the assessment is more proper because a hierarchical structure with too fine-grain subtopics will not help improve search ranking given the fact that there are only 10 ranking positions available on the first SERP.

### 4.2 Document Ranking Subtask

For the Document Ranking subtask, relevance judgment should be performed to result documents for all queries including clear, ambiguous and broad ones. To help assessors to finish the relevance judgment task, we extract all result documents in the pool from SogouT and ClueWeb. HTML documents are transformed into JPG version so that the appearance of documents to each assessor is the same. It can also reduce the efforts of assessors to load a Web page from its HTML version. The annotation process is completed in the following steps:

- Result pool construction: Due to limited annotation resources, we only cover a number of top results from a selection of submitted runs from participating groups. For the Chinese DR task, we choose top 20 results from runs with top priority from each group. While for English DR task, we choose The result top 10 results from runs with top priority from each group. The result pool for Chinese and English DR tasks contain 2,525 and 1,930 result documents, separately.
- Annotation task 1 (relevance judgment): for each document-query pair, the assessors are required to decide whether the document is relevant to the query with a 4-grade score (3: highly-relevant, 2: relevant, 1: irrelevant, 0: spam).
- Annotation task 2 (subtopic judgment) For a result document annotated as 2 or 3 in the first step for a broad or ambiguous query, the assessors should point out which first-level and second-level subtopic this document is relevant to. If one document isn't relevant to any of the subtopics, it shouldn't be regarded as a relevant one. For clear queries, there is no need to finish this step.



With the above procedure, we obtain the document relevance assessment result both to queries and to corresponding subtopics. For clear queries, the original NDCG score is calculated as the evaluation result. For ambiguous and broad queries, we choose corresponding first-level subtopics in the calculation of  $D\#$ -measures. Second-level subtopics are not involved in the evaluation of DR tasks because the number of subtopics are too many (about 50) for a practical Web search scenario.

### 4.3 TaskMine Subtask

//to be added by TaskMine organizers

## 5. OFFICIAL EVALUATION RESULTS

We will present the evaluation results in the following two sections. At first, Cranfield-like approach is adopted based on the result assessment described in Section 4. These results should be regarded as official results because they could be compared with existing testing results such as those in INTENT/INTENT2. The test collection could also be reused by researchers who do not participate in the IMine task. After that, we will show the user preference test results for Chinese DR task. Although those results could not be reused or compared with previous Cranfield-like evaluation results, we believe that comparison of these two results should help further our understanding in the research of diversified search evaluations.

### 5.1 Subtopic Mining Subtask

While reporting the evaluation results for the Subtopic Mining subtask, we will at first compare the performance of different participating groups in terms of  $Hscore$ ,  $Fscore$  and  $Sscore$ , separately. We will also test different parameters of H-measures. After that, we will show the evaluation results with H-measures for both ambiguous and broad queries.

#### 5.1.1 Hscore comparison

Comparison of  $Hscore$  of different participating runs are shown in Figure 1. According to the figure we can see that for the Chinese SM task, CNU performs best with  $Hscore$  of 0.5789. Meanwhile, best runs from THUSAM, FRDC and KLE also gain promising results. Significance test results (two-tailed t-Test with  $p$ -value<0.01) show that the best results of CNU, KLE, THUSAM and FRDC cannot be separated from each other. According to these participants' descriptions, clustering technique was adopted by most of these runs to group the provided candidates and word embedding as well as semantic expansion were also employed to extract subtopics.

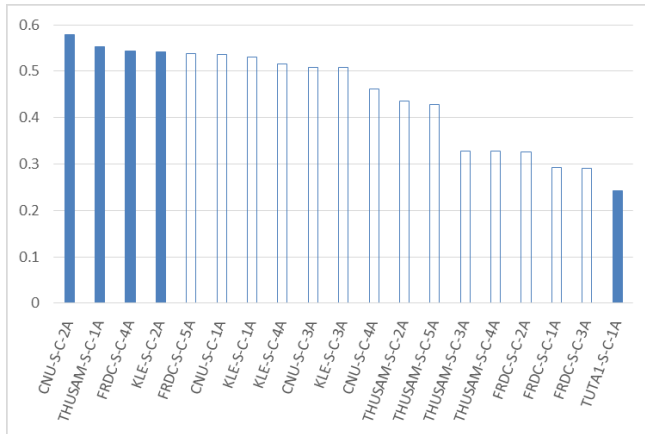


Figure 1.  $Hscores$  of submitted runs for unclear queries in

### Chinese Subtopic Mining (run with the highest performance for each participant is shown as a colored block while other runs are shown as non-colored blocks)

Figure 2 shows the  $Hscore$  distribution of proposed runs in English Subtopic Mining task. We can see that KUIDL and THUSAM gain best performances and their  $Hscores$  are much higher than those of other runs and their performance differences is not significant (two-tailed t-Test with  $p$ -value<0.01). According to their descriptions for submitted runs, KUIDL adopted the content from search engine result pages and THUSAM rely on Wikipedia page structures. One common feature from both runs is that first-level subtopics are always a sub-string for their corresponding second-level subtopics. The assessors tend to believe that this kind of second-level subtopics belong to the scope of first-level ones and annotate these subtopic pairs as correct ones.

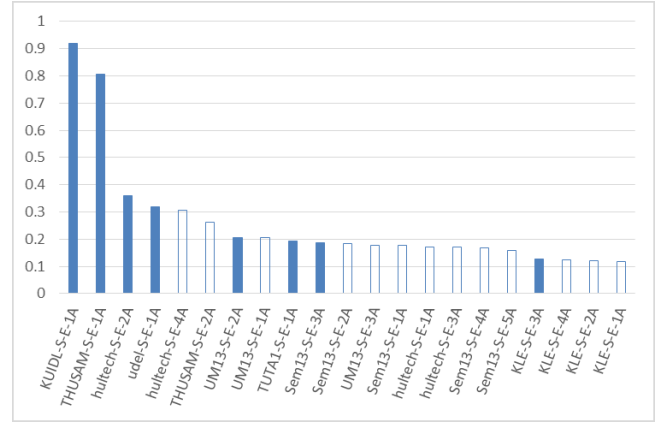


Figure 2.  $Hscores$  of submitted runs for unclear queries in English Subtopic Mining (run with the highest performance for each participant is shown as a colored block while other runs are shown as non-colored blocks)

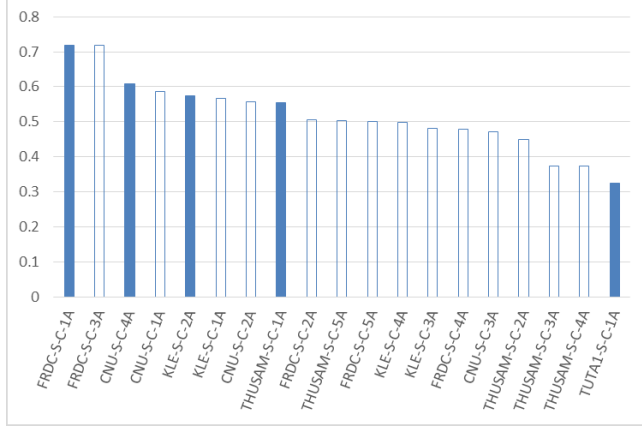
This is the first year that we introduce a hierarchical structure in subtopic extraction tasks. The relationship between first-level and second-level subtopics shares similar characteristics with the relationship between entities in knowledge graphs. Meanwhile, diversified search mainly focuses on covering more popular user interests behind these topics. From the above results, we can see that the best runs from Chinese SM task focus on clustering technique while those in English prefer candidate pairs in which first-level subtopics are substrings for corresponding second-level ones. We hope to see how the introduction of user behavior data (the organizers shared some user behavior data for Chinese SM task while participants can also acquire English/Japanese query frequency data from services such as google trends) could improve these methods in the future tasks or discussions.

#### 5.1.2 Fscore Comparison

$Fscore$  evaluates whether the submitted ranking lists of first-level subtopics meet users' diversified search intents. Comparison results for the participating runs are shown in Figures 3 and 4 for Chinese and English SM tasks. Note that only ambiguous queries are evaluated in this part because there may be several different groups of first-level subtopics that are all reasonable for broad queries.

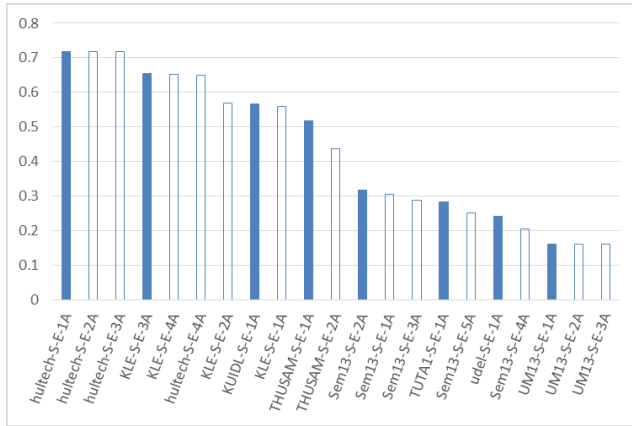
We can see that for Chinese SM task, FRDC gain highest  $Fscores$  with the runs FRDC-S-C-1A and FRDC-S-C-3A. Detailed analysis show that their runs gain both good I-recall (0.76 on average) and D-nDCG (0.67 on average) values. One interesting finding lies that their best performing run according to  $Hscore$  (FRDC-S-C-4A) fails to get high  $Fscore$  value while the two runs that gain best result

in *Fscore* (FRDC-S-C-1A and FRDC-S-C-3A) don't got promising results in *Hscores*, either. According to the system description provided by participant, we find that FRDC adopts the same strategy in developing first-level subtopics in FRDC-S-C-1A and FRDC-S-C-3A. Query subtopics provided by organizers as well as knowledge graph entries (from Baidu Baike) are adopted as candidates, which are clustered based on corresponding SERPs collected from Google. After that, new word detection techniques are employed to generate the first level subtopics.



**Figure 3. *Fscores* of submitted runs for ambiguous queries in Chinese Subtopic Mining (run with the highest performance for each participant is shown as a colored block while other runs are shown as non-colored blocks)**

For English SM task, we can see that hultech gains best performance in *Fscores*. While the difference between the best results from hultech, KLE and KUIDL are not significant (two-tailed t-Test with  $p\text{-value} < 0.01$ ). According to participants' descriptions, KLE and hultech both adopt pattern matching on the provided subtopic candidates to generate first-level subtopics. Meanwhile, KUIDL employs a different strategy by extracting first-level subtopics from top SERPs for given queries.



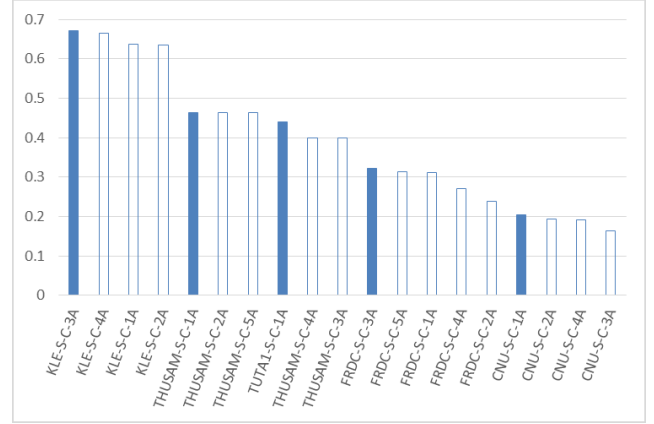
**Figure 4. *Fscores* of submitted runs for ambiguous queries in English Subtopic Mining (run with the highest performance for each participant is shown as a colored block while other runs are shown as non-colored blocks)**

### 5.1.3 *Sscore* Comparison

*Score* shows the fine-grained subtopic mining performance of participating runs. As stated in previous sections, at most 50 second-level subtopics are submitted in each run and they should be ranked within the whole query instead of within corresponding

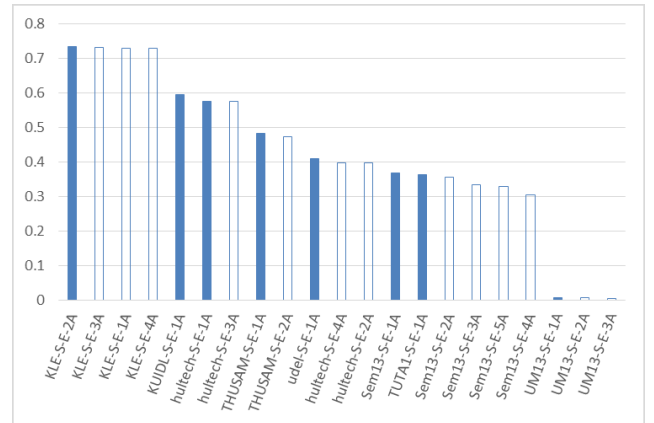
first-level subtopics. By this means, we could evaluate the system's performance in meeting fine-grained search intents.

According to the results shown in Figure 5, we can see that KLE obtains best *Sscore* performance in Chinese SM task. The difference between their best performing run and that from the second best group (THUSAM) is significant (two-tailed t-Test with  $p\text{-value} < 0.01$ ). From the descriptions provided by participants, we can see that the four runs submitted by KLE all adopt similar strategy (with different parameters). They are based on the provided subtopic candidates (query suggestion, query dimension, related queries and baseline documents) and combined with certain re-ranking techniques.



**Figure 5. *Sscores* of submitted runs for unclear queries in Chinese Subtopic Mining (run with the highest performance for each participant is shown as a colored block while other runs are shown as non-colored blocks)**

KLE also gains best performance in *Sscore* according to the English SM results shown in Figure 6. The difference between their best performing result and the second best one (from KUIDL) is also significant. We can see that a similar strategy (combination of candidates from different sources) is adopted in both Chinese and English mining tasks and their submitted runs are all based on this strategy with different parameters. We hope to see more details in participant's technical paper.



**Figure 6. *Sscores* of submitted runs for unclear queries in English Subtopic Mining (run with the highest performance for each participant is shown as a colored block while other runs are shown as non-colored blocks)**

### 5.1.4 *H-Measure* Comparison

With the *Hscore*, *Fscore* and *Sscore* result comparisons in previous

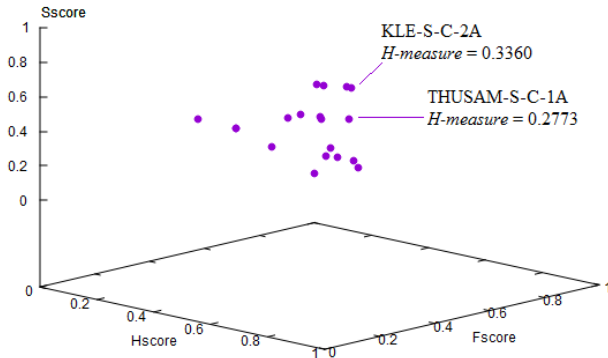


sections, we generate the *H-measure* results according to Equation (6) in Tables 7, 8 and 9. The best performing results are also shown in Figures 7, 8 and 9 so that we can see in what way these results outperform other runs. Note that the first-level subtopics for broad queries will not be taken into consideration in the evaluation because there may be various standards for organizing high-level aspects for these queries. For example, “harry potter movies/harry potter books/harry potter games” and “harry potter and the prisoner of Azkaban/harry potter and the goblet of fire/harry potter and the half blood prince” may both be good categories of subtopics for the query “harry potter” (IMINE 0083), but they lead to quite different first-level subtopic evaluation results. Therefore, for the “broad” queries in Table 6, the parameter  $\alpha$  in *H-measure* calculation is set to 0 while  $\beta$  is set to 1.0 in Equation (6). As stated in Section 3.1, we choose equal values of  $\alpha$  and  $\beta$  ( $\alpha = \beta = 0.5$ ) for ambiguous queries.

For Chinese SM task, KLE gain best performance with all four submitted runs. We can see that *Sscore* contributes most to their performance and they also gain nice results in *Hscores* and *Fscores*. As stated in Section 5.1.3, the four runs submitted by them all adopt similar strategy (with different parameters).

**Table 7. Chinese Subtopic Mining runs ranked by *H-measure* (official result) over 33 unclear topics. The highest value in each column is shown in bold.**

|               | <i>Hscore</i> | <i>Fscore</i> | <i>Sscore</i> | <i>H-measure</i> |
|---------------|---------------|---------------|---------------|------------------|
| KLE-S-C-2A    | 0.5413        | 0.5736        | 0.6339        | <b>0.3360</b>    |
| KLE-S-C-1A    | 0.5306        | 0.5666        | 0.6360        | 0.3303           |
| KLE-S-C-4A    | 0.5148        | 0.4986        | 0.6640        | 0.3279           |
| KLE-S-C-3A    | 0.5072        | 0.4817        | <b>0.6718</b> | 0.3255           |
| THUSAM-S-C-1A | 0.5527        | 0.5537        | 0.4634        | 0.2773           |
| THUSAM-S-C-5A | 0.4287        | 0.5040        | 0.4626        | 0.2224           |
| THUSAM-S-C-2A | 0.4347        | 0.4498        | 0.4633        | 0.2204           |
| FRDC-S-C-5A   | 0.5377        | 0.5004        | 0.3139        | 0.1757           |
| CNU-S-C-2A    | <b>0.5789</b> | 0.5569        | 0.1932        | 0.1748           |
| CNU-S-C-1A    | 0.5353        | 0.5867        | 0.2045        | 0.1739           |
| FRDC-S-C-4A   | 0.5436        | 0.4782        | 0.2715        | 0.1724           |
| CNU-S-C-4A    | 0.4611        | 0.6073        | 0.1910        | 0.1407           |
| THUSAM-S-C-4A | 0.3284        | 0.3744        | 0.3993        | 0.1404           |
| THUSAM-S-C-3A | 0.3284        | 0.3744        | 0.3981        | 0.1400           |
| FRDC-S-C-1A   | 0.2931        | <b>0.7191</b> | 0.3110        | 0.1327           |
| FRDC-S-C-3A   | 0.2897        | <b>0.7191</b> | 0.3214        | 0.1326           |
| CNU-S-C-3A    | 0.5086        | 0.4708        | 0.1626        | 0.1189           |
| TUTAI-S-C-1A  | 0.2419        | 0.3242        | 0.4391        | 0.1126           |
| FRDC-S-C-2A   | 0.3257        | 0.5045        | 0.2381        | 0.1032           |



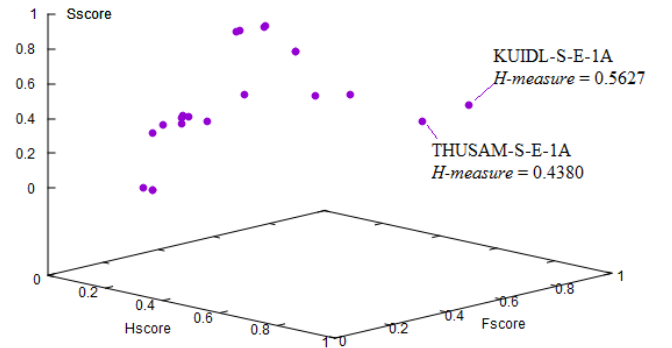
**Figure 7. Best performing runs in Chinese SM task and their performance comparison**

Different from the Chinese SM task, *Hscore* plays a central part in

English SM result comparisons. It is possibly due to the fact that there exist large differences between runs in *Hscore* for English SM task (see Figure 2). KUIDL-S-E-1A achieves both the best *Hscore* and best *H-measure* in Table 8. From Figures 4 and 6 we can also see that KUIDL’s *Fscore* and *Sscore* results are also quite nice compared with other runs. According to participant’s result descriptions, they try to extract hierarchical intents from search result landing pages’ structures. First-level subtopics are at first extracted from Web search results and then second-level ones are extracted by counting the co-occurrence of words in different page portions. It seems to be a rule-based method and we would like to see more details in participant’s technical paper.

**Table 8. English Subtopic Mining runs ranked by *H-measure* (official result) over 33 unclear topics. The highest value in each column is shown in bold.**

|                | <i>Hscore</i> | <i>Fscore</i> | <i>Sscore</i> | <i>H-measure</i> |
|----------------|---------------|---------------|---------------|------------------|
| KUIDL-S-E-1A   | <b>0.9190</b> | 0.5670        | 0.5964        | <b>0.5627</b>    |
| THUSAM-S-E-1A  | 0.8065        | 0.5179        | 0.4835        | 0.4380           |
| hultech-S-E-2A | 0.3596        | <b>0.7184</b> | 0.3977        | 0.1480           |
| hultech-S-E-4A | 0.3055        | 0.6496        | 0.3981        | 0.1323           |
| udel-S-E-1A    | 0.3202        | 0.2420        | 0.4103        | 0.1289           |
| THUSAM-S-E-2A  | 0.2634        | 0.4361        | 0.4732        | 0.1203           |
| KLE-S-E-3A     | 0.1273        | 0.6539        | 0.7317        | 0.1000           |
| KLE-S-E-4A     | 0.1242        | 0.6511        | 0.7294        | 0.0959           |
| KLE-S-E-2A     | 0.1194        | 0.5698        | <b>0.7342</b> | 0.0921           |
| KLE-S-E-1A     | 0.1179        | 0.5591        | 0.7298        | 0.0914           |
| hultech-S-E-1A | 0.1703        | <b>0.7184</b> | 0.5754        | 0.0888           |
| hultech-S-E-3A | 0.1703        | <b>0.7184</b> | 0.5754        | 0.0888           |
| TUTAI-S-E-1A   | 0.1933        | 0.2833        | 0.3647        | 0.0688           |
| Sem13-S-E-1A   | 0.1762        | 0.3043        | 0.3689        | 0.0634           |
| Sem13-S-E-2A   | 0.1844        | 0.3174        | 0.3566        | 0.0610           |
| Sem13-S-E-3A   | 0.1860        | 0.2882        | 0.3333        | 0.0606           |
| Sem13-S-E-4A   | 0.1672        | 0.2056        | 0.3039        | 0.0501           |
| Sem13-S-E-5A   | 0.1580        | 0.2511        | 0.3285        | 0.0470           |
| UM13-S-E-2A    | 0.2064        | 0.1624        | 0.0059        | 0.0016           |
| UM13-S-E-1A    | 0.2056        | 0.1624        | 0.0059        | 0.0015           |
| UM13-S-E-3A    | 0.1766        | 0.1624        | 0.0049        | 0.0014           |



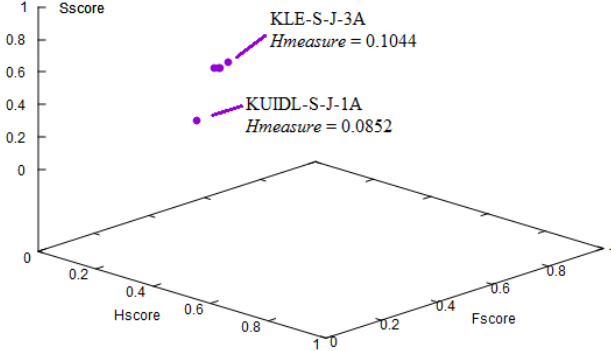
**Figure 8. Best performing runs in English SM task and their performance comparison**

As for Japanese SM task, since there are only two participating groups, we just compare their *H-measure* performances in this section and don’t present the *Hscore*, *Fscore* and *Sscore* comparisons, separately. From the results shown in Table 9 and Figure 9 we can see that KLE gain better performance than KUIDL in *H-measure* but the difference between their best performing runs is not significant (two-tailed t-Test with  $p$ -value<0.01). From the descriptions in submitted runs, we find that both KLE and KUIDL adopt similar strategies in different languages. They gain best

performance in Chinese SM and English SM, separately. We expect the two participating groups to compare their runs in different languages and hope to see some interesting findings.

**Table 9. Japanese Subtopic Mining runs ranked by  $H$ -measure (official result) over 34 unclear topics. The highest value in each column is shown in bold.**

|              | $H$ score     | $F$ score     | $S$ score     | $H$ -measure  |
|--------------|---------------|---------------|---------------|---------------|
| KLE-S-J-3A   | 0.2030        | <b>0.4745</b> | <b>0.5086</b> | <b>0.1044</b> |
| KLE-S-J-4A   | 0.2025        | 0.4248        | 0.4997        | 0.1014        |
| KLE-S-J-2A   | 0.1867        | 0.4609        | 0.4697        | 0.0907        |
| KLE-S-J-1A   | 0.1794        | 0.4662        | 0.4625        | 0.0864        |
| KUIDL-S-J-1A | <b>0.2702</b> | 0.2883        | 0.2848        | 0.0852        |



**Figure 9. Submitted runs in Japanese SM task and their performance comparison**

## 5.2 Document Ranking Subtask

As stated in Section 3.1, we follow the settings in INTENT/INTENT2 and choose  $D\#nDCG$  as the main evaluation metric for Document Ranking subtask. Since a hierarchy of subtopics is provided for each unclear query topic, we actually have two lists of subtopics for each of these queries: a first-level subtopic list and a second-level one. Therefore, we could evaluate the submitted runs with either fine-grained or coarse-grained search intents. The evaluation results are shown in Tables 10 and 11 for Chinese DR and English DR tasks, separately. In the results shown in these tables, the performance of clear queries are evaluated with  $nDCG$ , which could be regarded as a special case for  $D\#nDCG$  with no diversified subtopic lists.

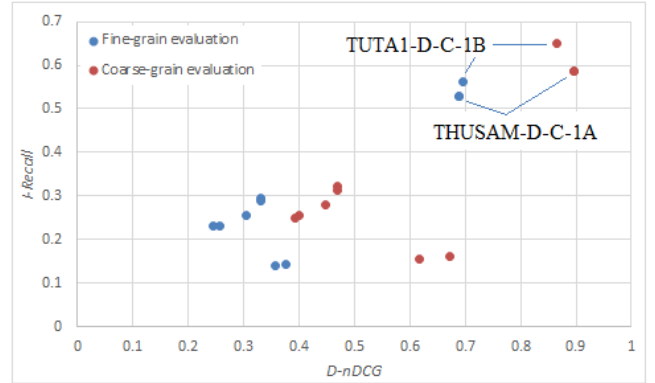
**Table 10. Chinese Document Ranking runs ranked by coarse-grain  $D\#nDCG$  (official result) over all query topics. The highest value in each column is shown in bold.**

|               | Coarse-grain results<br>(evaluated with<br>first-level subtopics) | Fine-grain results<br>(evaluated with<br>second-level subtopics) |
|---------------|---|--|
| TUTA1-D-C-1B  | <b>0.7334</b>   | <b>0.6538</b>  |
| THUSAM-D-C-1A | 0.6965  | 0.6127   |
| THUSAM-D-C-1B | 0.6943  | 0.6106   |
| FRDC-D-C-1A   | 0.4619  | 0.4118   |
| FRDC-D-C-3A   | 0.4440  | 0.3950   |
| FRDC-D-C-2A   | 0.3899  | 0.3402   |
| FRDC-D-C-5A   | 0.3841  | 0.3338   |
| FRDC-D-C-4A   | 0.3746  | 0.3240   |
| THUSAM-D-C-2B | 0.3697  | 0.2711   |
| THUSAM-D-C-2A | 0.3502  | 0.2623   |

Evaluation results in Tables 10 and 11 show that the coarse-grain results and fine-grain results are highly correlated (correlation values are both over 0.99). In Chinese DR task, TUTA gains best performance for both coarse-grain and fine-grain subtopic lists and

the difference between their best run (TUTA1-D-C-1B) and the second best run (THUSAM-D-C-1A) is significant (two-tailed t-Test with  $p$ -value<0.01). According to descriptions given by TUTA, they adopt the subtopic list submitted to Chinese SM task and use different ranking strategies for different kinds of topics. This run is based on the non-diversified baseline provided by organizers. Considering the fact that TUTA doesn't gain very promising results in SM task (no better than FRDC and THUSAM), we believe that the ranking strategy they adopt must be effective and we would like to read more details in the technical paper.

From the results in Figure 10, we can see that the coarse-grain  $D\#nDCG$  value of THUSAM-D-C-1A is higher than that of TUTA-D-C-1B. It probably show that the THUSAM run tends to adopt a relevance-oriented strategy while the TUTA one focuses more on intent recall. We can also see that these two runs gain much better performance than the other runs according to both coarse-grain and fine-grain results.



**Figure 10. Best performing runs in Chinese DR task and their relationship with other submitted runs**

According to evaluation results in Table 11, udel gains best performance with coarse-grain subtopic lists but the differences among best runs of udel, UM13, TUTA1 and Sem13 are not significant (two-tailed t-Test with  $p$ -value<0.01). Similarly,  $D\#nDCG$ s of the best performing runs of udel, UM13, TUTA1 and Sem13 with fine-grain subtopic lists are not significantly different, either. It is probably due to the fact that the pool depth for English DR runs are a bit shallow (covers top 10 results of the top priority runs). For the top performing runs, udel adopts query suggestions as inputs and use data fusion techniques to combine different ranking lists. TUTA adopts the same strategy in Chinese DR and UM13 employs a number of external resources including query logs, Wikipedia, ConceptNet and query suggestions from commercial search engines.

**Table 11. English Document Ranking runs ranked by coarse-grain  $D\#nDCG$  (official result) over all query topics. The highest value in each column is shown in bold.**

|              | Coarse-grain results<br>(evaluated with first-<br>level subtopics) | Fine-grain results<br>(evaluated with second-<br>level subtopics) |
|--------------|--|---|
| udel-D-E-1A  | <b>0.6297</b>  | 0.5469  |
| UM13-D-E-1A  | 0.6254   | 0.5566  |
| TUTA1-D-E-1B | 0.6170   | <b>0.5668</b>   |
| Sem13-D-E-1A | 0.6022   | 0.5291  |
| UM13-D-E-2A  | 0.6001   | 0.5309  |
| Sem13-D-E-3A | 0.4735   | 0.3985  |
| Sem13-D-E-2A | 0.4495   | 0.3806  |
| UM13-D-E-3A  | 0.4474   | 0.3770  |

|              |        |        |
|--------------|--------|--------|
| udel-D-E-2A  | 0.3900 | 0.3181 |
| udel-D-E-4A  | 0.3472 | 0.2808 |
| TUTA1-D-E-2B | 0.3314 | 0.2601 |
| Sem13-D-E-4A | 0.3227 | 0.2505 |
| Sem13-D-E-5A | 0.3081 | 0.2414 |
| udel-D-E-3A  | 0.0985 | 0.0784 |
| udel-D-E-5A  | 0.0932 | 0.0877 |

From the results shown in Figure 11, we can see that udel-D-E-1A gains much higher I-recall value compared with other runs in coarse-grain evaluation. It also gains best D-nDCG value in fine-grain evaluation. This difference shows that although the results given by coarse-grain and fine-grain evaluations are highly correlated, same strategy results in different performance evaluation results with subtopic lists in different grains.

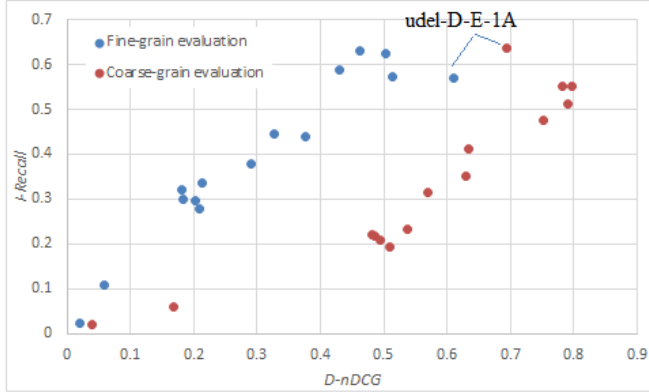


Figure 11. Best performing runs in Chinese DR task and their relationship with other submitted runs

### 5.3 TaskMine Subtask

//to be added by TaskMine organizers

## 6. USER PREFERENCE TEST RESULTS

//to be added in the final version

## 7. CONCLUSIONS AND FUTURE WORK

IMine task aims to mine users' diversified intents behind their simple, unspecified and sometimes ambiguous queries submitted to search engines. It follows the research framework of previous INTENT/INTENT2 tasks in NTCIR9/10 but features the mining of hierarchical subtopic structures. In this year's task, the organizers work with participants to develop new evaluation metrics for SM task (*Hscore*, *Fscore*, *Sscore* and *H-measure*) and employ them to evaluate the performance of this Web search intent mining task. The evaluation of DR task is also different from previous tasks in that both a fine-grain and a coarse-grain comparison can be obtained with the generated second-level and first-level subtopic lists, separately.

Through the evaluation results, we found that best performing runs for SM task in different languages are usually based on a combination of different information resources (query suggestions, query dimensions and related queries). *Hscore* plays a central role in English SM task but the best runs in Chinese and Japanese runs seem to be *Sscore*-oriented. We hope to obtain more interesting findings by comparison of different participants' methods after the participants' technical papers are available.

Although the pool depth for DR task is not so deep and the

evaluation results are not significantly different from each other for English task, we still find several interesting findings through the evaluation process. We find that coarse-grain results and fine-grain results are highly correlated, which means that it may not be necessary to use fine-grain subtopic list to avoid extra annotation efforts.

## 8. ACKNOWLEDGMENTS

We appreciate Prof. Jamie Callan and his team for providing the ClueWeb collection, which dramatically reduces the working efforts of participants in Document Ranking subtask. We also benefit a lot from discussions with Dr. Tetsuya Sakai and the PC chairs of NTCIR-11.

## 9. REFERENCES

- [1] R. Song, M. Zhang, T. Sakai, M. P. Kato, Y. Liu, M. Sugimoto, Q. Wang, and N. Orii. Overview of the NTCIR-9 INTENT Task. In Proceedings of NTCIR-9. 2011, 82-105.
- [2] T. Sakai, Z. Dou, T. Yamamoto, Y. Liu, M. Zhang, and R. Song. Overview of the NTCIR-10 INTENT-2 Task. In Proceedings of NTCIR-10 Workshop Meeting. 2013.
- [3] Zhicheng Dou, Sha Hu, Yulong Luo, Ruihua Song and Ji-Rong Wen.: Finding Dimensions for Queries, ACM CIKM 2011, October 2011.
- [4] Yiqun Liu, Junwei Miao, Min Zhang, Shaoping Ma, Liyun Ru.: How Do Users Describe Their Information Need: Query Recommendation based on Snippet Click Model. Expert Systems With Applications. 38(11): 13847-13856, 2011.
- [5] M. Kendall. A new measure of rank correlation. Biometrika, 30:81-89, 1938.
- [6] T. Sakai. Evaluating evaluation metrics based on the bootstrap. In In Proceedings of ACM SIGIR 2006. ACM, Seattle, Washington, USA, pages 525-532, August 2006.
- [7] T. Sakai. Evaluation with informational and navigational intents. In In Proceedings of ACM WWW 2012. ACM, Lyon, France, pages 499-508, April 2012.
- [8] Ruihua Song, Zhenxiao Luo, Jian-Yun Nie, Yong Yu, and Hsiao-Wuen Hon. Identification of ambiguous queries in web search. Information Processing & Management, 45(2):216-229, 2009.
- [9] Sakai, T. and Song, R.: Evaluating Diversified Search Results Using Per-Intent Graded Relevance, ACM SIGIR 2011, pp.1043-1052, July 2011.
- [10] R. Agrawal, S. Gollapudi, A. Halverson, and S. Ieong. (2009). Diversifying Search Results. In Proceedings of the Second ACM International Conference on Web Search and Data Mining. ACM, Barcelona, Spain (pp.5-14).
- [11] O. Chapelle, D. Metzler, Y. Zhang, and P. Grinspan. (2009). Expected Reciprocal Rank for Graded Relevance. In Proceedings of ACM CIKM 2009. ACM, Hong Kong, China. 621-630.
- [12] T. Sakai, N. Craswell, R. Song, S. Robertson, Z. Dou, and C.-Y. Lin. (2010). Simple Evaluation Metrics for Diversified Search Results. In 3rd International Workshop on Evaluating Information Access. Tokyo, Japan (pp.42-50).
- [13] Cheng Zhai, William W. Cohen, John Lafferty. (2003). Beyond Independent Relevance: Methods and Evaluation Metrics for Subtopic Retrieval. In Proceedings of ACM SIGIR 2003. ACM, Toronto, Canada (pp.10-17).

- [14] T. Sakai. (2012). Evaluation with Informational and Navigational Intents. In Proceedings of ACM WWW 2012. ACM, Lyon, France (pp.499-508).