

# Search Result Diversity Evaluation based on Intent Hierarchies

Xiaojie Wang, Ji-Rong Wen, Zhicheng Dou, Tetsuya Sakai, and Rui Zhang.

**Abstract**—Search result diversification aims at returning diversified document lists to cover different user intents of a query. Existing diversity measures assume that the intents of a query are disjoint, and do not consider their relationships. In this paper, we introduce intent hierarchies to model the relationships between intents, and present four weighing schemes. Based on intent hierarchies, we propose several hierarchical measures that take into account the relationships between intents. We demonstrate the feasibility of hierarchical measures by using a new test collection based on TREC Web Track 2009-2013 diversity test collections and by using NTCIR-11 IMine test collection. Our main experimental findings are: (1) Hierarchical measures are more discriminative and intuitive than existing measures. In terms of intuitiveness, it is preferable for hierarchical measures to use the whole intent hierarchies than to use only the leaf nodes; (2) The types of intent hierarchies used affect the discriminative power and intuitiveness of hierarchical measures. We suggest the best type of intent hierarchies to be used according to whether the nonuniform weights are available; (3) To measure the benefits of the diversification algorithms which use automatically mined hierarchical intents, it is important to use hierarchical measures instead of existing measures.

**Index Terms**—Ambiguity, Diversity, Evaluation, Novelty, Hierarchy.

## 1 INTRODUCTION

Nowadays, People tend to meet their daily information needs by issuing keywords into search engines. However, these keywords, i.e. queries, are often ambiguous or broad [1], [2], [3], [4]. The queries usually have several interpretations or aspects, also known as subtopics or user intents. When users submit the same query to retrieval systems, they may want different information returned to fulfill their own information needs. This poses a challenge to search engines when the user intent cannot be known in advance.

To tackle this problem, a wide range of search result diversification algorithms ([5], [6], [7], [8], [9], [10], [11], [12], [13], [14]) have been proposed over the past years. They aim at returning a diversified ranked document list that covers different intents of a query. In the meantime, some researchers have introduced a variety of diversity measures, such as I-rec [15],  $\alpha$ -nDCG [16], Intent-Aware measures [5], D $\sharp$ -measures [17], etc. These measures evaluate ranked lists in terms of both diversity and relevance, and can be used to indicate which diversification algorithms are better. Existing diversity measures assume that the users' information needs could be represented by a single layer of intents and different types of intents are independent of each other.

However, intents can be related to each other in reality, which is illustrated as follows.

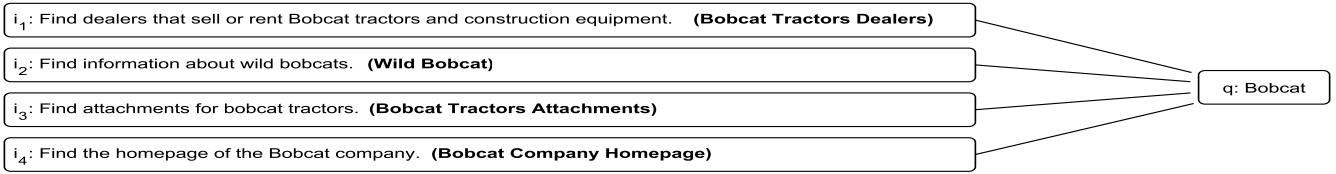
We use the query “bobcat”, No. 77 topic in Text Retrieval Conference(TREC) 2010 Web Track [18], as an example. This query is ambiguous because of the polysemy of “bobcat”: one interpretation is a company called “bobcat company” whose core business is about tractors; another interpretation is a kind of wild animals called “wild bobcat.” We show its official intents, marked by  $i_1$ - $i_4$ , in Figure 1(a). The figure shows that except intent  $i_2$  which is about “wild bobcat,” the remaining ones,  $i_1$ ,  $i_3$ , and  $i_4$ , are all about “bobcat company.” This indicates that  $i_1$ ,  $i_3$ , and  $i_4$  are more related to each other, but are less related to  $i_2$ . Even within the three intents about “bobcat company,”  $i_1$  and  $i_3$  are closer because they are about the businesses involving tractors of the company, whereas  $i_4$  is about homepage the company. We argue that this kind of relationships between intents should be modeled when evaluating search result diversity. However, none of existing measures considers this.

Specifically, we find two submitted runs for the query, cmuFuTop10D and THUIR10DvNov, in TREC Web Track 2010 diversity task. cmuFuTop10D covers  $i_1$ ,  $i_3$ , and  $i_4$ , while THUIR10DvNov covers  $i_1$ ,  $i_2$ , and  $i_4$  in their top ten documents. Since  $i_1$ ,  $i_3$ , and  $i_4$  are all about “bobcat company,” cmuFuTop10D misses another interpretation of bobcat, i.e. “wild bobcat,” but THUIR10DvNov covers both interpretations. In this sense, the latter is more diversified but I-rec [15] treats them as equally good because they cover the same number of intents. Some other existing measures also have similar problems, which will be illustrated in Section 3.4.2. We think that this is due to their lack of recognition of the relationships among intents.

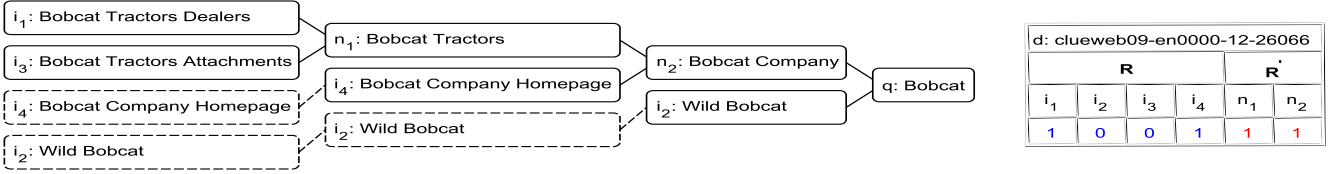
In light of the above observation, we introduce intent hierarchies to model the relationships among intents. We design hierarchical measures using the intent hierarchies

- Xiaojie Wang is with the School of Information at Renmin University of China, and the School of Computing and Information Systems at The University of Melbourne. E-mail: xiaojiew1@student.unimelb.edu.au
- Zhicheng Dou and Ji-Rong Wen are with the School of Information, Beijing Key Laboratory of Big Data Management and Analysis Methods, and DEKE, Renmin University of China, Beijing 100872, P.R. China. E-mail: dou@ruc.edu.cn, jirong.wen@gmail.com
- Tetsuya Sakai is with Department of Computer Science and Engineering, Waseda University. E-mail: tetsuyasakai@acm.org
- Rui Zhang is with the School of Computing and Information Systems at The University of Melbourne. E-mail: rui.zhang@unimelb.edu.au

Manuscript received August 19 2016; revised July 14, 2017.



(a) Official intents of the query “bobcat”.



(b) LEFT: OIH is comprised of the solid boxes, whereas EIH includes both solid and dashed nodes. RIGHT: An example showing relevance assessments for the added nodes (under  $R'$  in red) derived from relevance assessments for the official intents (under  $R$  in blue).

Fig. 1. The official intents, original intent hierarchy (OIH), and extended intent hierarchy (EIH) of No. 77 query “bobcat” in TREC Web Track 2010.

to solve the problems mentioned above. This paper is the extended version of the SIGIR 2016 paper [19]. In the original work [19], we have found that hierarchical measures, especially those using the whole intent hierarchy, are better than existing measures, which uses an intent list, in terms of discriminative power and intuitiveness. The main extensions of this journal version are:

(1) We propose four weighting schemes that are used to model the node weights in an intent hierarchy. We examine the impact of the different types of intent hierarchies, including whether the leaf nodes have the same depth and which weighting scheme is used, for hierarchical measures in terms of discriminative power and intuitiveness. In the experiments, we find the best type of intent hierarchies when nonuniform weights are available and when only uniform weights are known respectively.

(2) We find that it is crucial for hierarchical diversification algorithms to be evaluated by hierarchical measures. The benefits in search result diversification by re-ranking the results to cover the automatically generated hierarchical intents as much as possible may be invisible to existing measures that measure the diversity using intent lists. Hierarchical diversification algorithms show more gains when evaluated by hierarchical measures than existing measures.

(3) Besides TREC Web Track 2009-2013 test collections, we also experiment with NTCIR-11 IMine<sup>1</sup> test collection that has official intent hierarchies with nonuniform weights. We find that the conclusions drawn from the two sources of test collections are consistent with each other.

(4) We reveal that Layer-Aware measures may be counterintuitive because of their preference of high relevance to popular nodes. The experiments confirm that they are less intuitive among the proposed hierarchical measures.

The remainder of this paper is organized as follows. Section 2 describes some existing diversity measures and the methods for testing them. In Section 3, we introduce intent hierarchies, and our method for creating a new test collection based on TREC Web Track test collections. We then propose several new diversity measures that can use

the intent hierarchies. Section 4 describes the experimental results and analysis. We conclude our work in Section 5.

## 2 RELATED WORK

Evaluation measures play an important role in the scientific research because they serve as the inexpensive methods for monitoring the technological progress. In the Information Retrieval experiments, evaluation measures use test collections to evaluate system performances. Depending on the task at hand, it is essential to analyze the properties of evaluation measures and use the appropriate ones. Search result diversification aims to cover different intents by a ranked list. Given a query  $q$ , most existing measures evaluate the diversified search results by modeling users’ information needs as a flat list of intents  $fig$ . Some measures can handle intent probability  $Pr(ijq)$  and graded relevance assessments but some cannot. In this section, we summarize the previous work on designing and evaluating diversity measures.

### 2.1 Diversity Measures

#### 2.1.1 Intent Recall

Intent recall(I-rec) [15], also known as subtopic recall [20] is the proportion of intents covered by a ranking list. Let  $d_r$  denote the document at rank  $r$ , and let  $I(d_r)$  denote the set of intents to which document  $d_r$  is relevant. Then,  $I-rec$  for a certain cutoff  $K$  can be expressed as:

$$I-rec@K = \frac{\sum_{r=1}^K I(d_r)j}{jfig} \quad (1)$$

The idea of I-rec is to credit minor intents. I-rec does not consider the positions of relevant documents, and cannot handle intent probability and graded relevance assessments.

#### 2.1.2 $\alpha$ -nDCG

In order to balance both relevance and diversity of ranked lists,  $\alpha$ -nDCG [16] is defined as:

$$-nDCG@K = \frac{\sum_{r=1}^K NG(r) = \log(r+1)}{\sum_{r=1}^K NG(r) = \log(r+1)} \quad (2)$$

$$NG(r) = \frac{J_i(r)(1 - \alpha)^{C_i(r-1)}}{i2fig}$$

1. <http://www.thuir.org/IMine/>

where  $NG(r)$  is  $NG(r)$  in the ideal ranked list;  $J_i(r)$  is 1 if the document at rank  $r$  is relevant to intent  $i$ , and 0 otherwise;  $C_i(r) = \sum_{k=1}^r J_i(k)$  is the number of relevant documents to intent  $i$  within top  $r$ ; and  $\alpha$  is a parameter.

### 2.1.3 Intent-Aware measures

Intent-Aware measures (IA measures) [5] is a general framework to evaluate ranked document lists. Assuming that  $\prod_{i \in \text{fig}} Pr(ijq) = 1$ ,  $M$ -IA can be computed as:

$$M\text{-IA}@K = \prod_{i \in \text{fig}} Pr(ijq) M_i@K \quad (3)$$

where  $M_i$  is the per-intent version of measure  $M$ . Measure  $M$  can be nDCG [21], ERR [22], nERR [23], etc.

### 2.1.4 D-measures

Assume that  $g_i(r)$  is the gain value of the document at rank  $r$  for intent  $i$ , and  $g_i(r)$  is calculated using per-intent relevance assessments. Then the global gain at rank  $r$  is given by:

$$GG(r) = \prod_{i \in \text{fig}} Pr(ijq) g_i(r) \quad (4)$$

Let  $CGG(r) = \prod_{k=1}^r GG(k)$ , i.e. the cumulative global gain at rank  $r$ . Let  $GG(r)$  and  $CGG(r)$  denote the global gain and the cumulative global gain at rank  $r$  in the ideal ranked list. The ideal list is obtained by listing up all relevant documents in descending order of global gains. Let  $J(r) = 1$  if the document at rank  $r$  is relevant to any of the intents  $\text{fig}$ , and  $J(r) = 0$  otherwise. Let  $C(r) = \sum_{k=1}^r J(k)$ , which is the number of relevant documents within top  $r$ .  $D$ -nDCG and  $D$ -Q at document cutoff  $K$  are defined as:

$$D\text{-nDCG}@K = \frac{\prod_{r=1}^K GG(r) = \log(r+1)}{\prod_{r=1}^K GG(r) = \log(r+1)} \quad (5)$$

$$D\text{-Q}@K = \frac{1}{\min(K; R)} \prod_{r=1}^K J(r) \frac{C(r) + CGG(r)}{r + CGG(r)} \quad (6)$$

where  $R$  is the number of judged relevant documents.

### 2.1.5 D<sub>#</sub>-measures

D<sub>#</sub>-measures [17] aim to boost intent recall, and to reward documents that are highly relevant to popular intents.  $D_{\#}$ -measure is defined as:

$$D_{\#}\text{-measure}@K = I\text{-rec}@K + (1 - \gamma) D\text{-measure}@K \quad (7)$$

where  $D$ -measure can be D-nDCG or D-Q.  $\gamma$  is a parameter between 0 and 1. D<sub>#</sub>-measures are free of the under-normalization problem of  $\alpha$ -nDCG and IA measures.

The measures mentioned above are widely used in several tasks of TREC Web Track <sup>2</sup> or NII Testbeds and Community for Information access Research (NTCIR) <sup>3</sup>, but they do not take the relationships between intents into consideration, which is what we aim to deal with in this paper.

## 2.2 Evaluation of Diversity Measures

Given a significance level, *discriminative power* measures the stability of measures across queries and experiments based on significance tests, e.g. paired bootstrap test [24], Tukey's Honestly Significant Differences(HSD) [25] test, etc.

*Concordance test* [26] is proposed to quantify the intuitiveness of diversity measures. In concordance test, one or more gold standard measures are chosen and assumed to truly represent intuitiveness. Given two diversity measures  $M_1$  and  $M_2$ , the relative intuitiveness of  $M_1$  (or  $M_2$ ) is measured in terms of preference agreement with the gold standard measures. The preference agreement is that  $M_1$  (or  $M_2$ ) agrees with the gold standard measure(s) about which one of two ranked lists should be preferred.

Kendall's  $\tau$  [27] is a statistic to measure the *rank correlation* of two rankings. However,  $\tau$  lacks the property of top heaviness. In the context of IR evaluation, the swaps near the top is generally more important than those near the bottom.  $\tau_{ap}$  [28] is proposed to deal with the problem. Note that  $\tau$  is symmetric but  $\tau_{ap}$  is not. However, a symmetric  $\tau_{ap}$  can be obtained by averaging two  $\tau_{ap}$  values when each list is treated as the former one. Both  $\tau$  and  $\tau_{ap}$  range from -1, which implies two ranked lists perfectly disagree, to 1, which implies two ranked lists are identical.

## 3 PROPOSED METHODS

In this section, we define two kinds of intent hierarchies with four weighting schemes to represent the relationships between user intents. We then introduce our method for creating such intent hierarchies and obtaining the relevance assessments based on TREC Web Track 2009-2013 diversity test collections. Last, we propose several diversity measures based on intent hierarchies, and show that the new measures outperform the corresponding existing measures.

### 3.1 Intent Hierarchies

Given a query  $q$ , the users' information needs are represented as a set of intents  $\text{fig}$ . We assume these intents cannot be further subdivided, and refer to them as *atomic intents*. Based on the semantic relatedness of the intents, we build an intent hierarchy possessing the following properties:

**Property 1.** The intent hierarchy is in a tree structure, where every child has only one parent.

**Property 2.** The root of intent hierarchy is denoted by  $q$  itself, which stands for the users' information needs as a whole. The root is a dummy node for the completeness of the tree, and is not considered in our measures.

**Property 3.** When  $q$  is broad, the intent hierarchy is built in such a way that a parent node refers to a more general concept than its children, and a child node refers to one aspect of its parent. When  $q$  is ambiguous, each child node of the root denotes one interpretation of  $q$ , and each of its subtrees is built in the same way as a broad query.

**Property 4.** These atomic intents, i.e.  $\text{fig}$ , correspond one to one with leaves of the intent hierarchy. This means the number of leaves in the intent hierarchy is the same as the number of the atomic intents.

2. <http://plg.uwaterloo.ca/trecweb/>

3. <http://research.nii.ac.jp/ntcir/index-en.html>

An intent hierarchy that satisfies the above four properties is called an *original intent hierarchy (OIH)*. The leaf nodes of an OIH may not have the same depth. We extend such an OIH by recursively adding one child node to the leaves which do not have the highest depth until it satisfies:

**Property 5.** All leaf nodes of the intent hierarchy, i.e. the atomic intents, have the same depth.

The resulting intent hierarchy are called an *extended intent hierarchy (EIH)*. As an OIH, the leaves of an EIH also correspond one to one with the atomic intents *fig.*

We define an atomic intent subset  $S$  as a subset of atomic intents which has at least two atomic intents. Extending OIH to EIH is justified as follows: (1) An atomic intent subset  $S$  corresponds to  $|S|$  leaves in an OIH or EIH. The maximum depth of the nodes that are common ancestors of these leaves acts as a simple indicator of how much these intents are related to each other. Hence, we can refer to the maximum depth as the *redundant degree* of  $S$ ; (2) When OIH is extended to EIH, the redundant degree remains the same for any atomic intent subset. This means that if OIH, which is built by experts, thinks that one atomic intent subset is more related than another, so does EIH and vice versa.

We consider the root of an intent hierarchy as the zeroth layer, the child nodes of the root as the first layer and so forth. If an intent hierarchy only has the zeroth layer and the first layer, the height of the intent hierarchy is one. In the paper, a *single-layer intent hierarchy* refers to an intent hierarchy whose height is one, while a *multilayer intent hierarchy* refers to that whose height is greater than one.

### 3.2 Weighting Intent Hierarchies

For an intent hierarchy, its nodes are weighted according to their relative popularity. The node weights should satisfy:

**Property 6.**  $W(q) = 1$  and  $W(n) = \frac{1}{|C(n)|} \sum_{c \in C(n)} W(c)$  ( $\forall n \in N$ )

where  $q$  is the root node,  $N$  is the set of non-terminal nodes, and  $C(n)$  is the set of child nodes of node  $n$ .

We propose four weighing schemes that make resulting node weights satisfy Property 6 as follows:

(1) *Uniformly top-down (UT)* can be used in any situation. It assumes that given a parent node, its child nodes are equally weighted. Starting from the root, the weight of node  $n$  is  $W(n) = W(p)/|C(p)|$ , where  $p$  is the parent node of  $n$ , and  $C(p)$  is the set of child nodes of  $p$ .

(2) *Uniformly bottom-up (UB)* can be used in any situation. It assumes that given an intent hierarchy, its leaf nodes are equally weighted, i.e.  $W(n) = \frac{1}{|B|}$  ( $\forall n \in B$ ) where  $B$  is the set of leaf nodes. Starting from the leaves, the weight of node  $n$  is defined as:  $W(n) = \frac{1}{|C(n)|} \sum_{c \in C(n)} W(c)$ , where  $C(n)$  is the set of child nodes of node  $n$ .

(3) *Nonuniformly top-down (NT)* can only be used when all the node weights in an intent hierarchy are known but they do not satisfy Property 6. This may happen when the atomic intents that receive no relevant documents are removed and node weights need renormalization. Starting from the root, the new weight of node  $n$  is defined as:

$$W(n) = \frac{W^o(n) \cdot W(p)}{\sum_{c \in C(p)} W^o(c)} \quad (8)$$

where  $p$  is the parent of  $n$ ,  $C(p)$  is the children of  $p$ ,  $W^o(n)$  ( $W^o(c)$ ) is the original weight of node  $n$  ( $c$ ).

(4) *Nonuniformly bottom-up (NB)* can be used if only the weights of leaf nodes in an intent hierarchy are known. First the weights of leaf nodes are normalized, and then starting from the leaves, the weight of node  $n$  is  $W(n) = \frac{1}{|C(n)|} \sum_{c \in C(n)} W(c)$ , where  $C(n)$  is the set of child nodes of  $n$ .

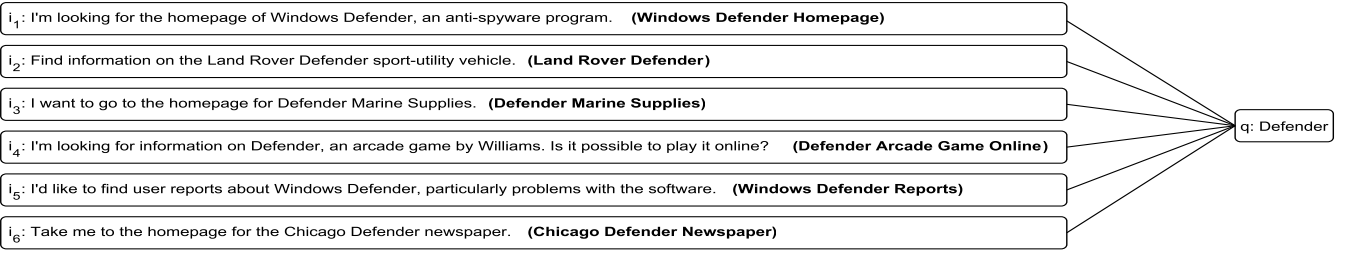
### 3.3 Creating Intent Hierarchies

Our proposed hierarchical measures, which will be introduced in 3.4, can work with any type of intent hierarchies. NTCIR-11 IMine test collection comes with multilayer intent hierarchies, and we directly use them in the experiments. In this section, we show the method for creating multilayer intent hierarchies from the predefined intents on TREC Web Track 2009-2013 diversity test collections. For each query in the test collections, the first intent is the same as the description of the query itself. Although the descriptions are the same, if a query has several different interpretations, the first intent is just one of these interpretations. A query's first intent does not refer to a more general concept than the other intents. So we do not treat the first intent differently.

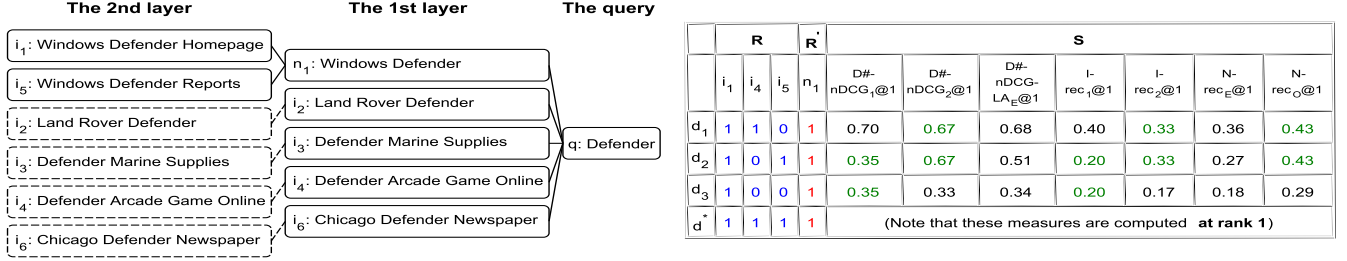
We use the official intents as atomic intents to avoid reassessing the relevance of documents. First we create OIH by manually grouping the official intents based on their semantic similarity. Then, we extend them to EIH. Figure 1 illustrates how we create OIH and EIH for the query "bobcat" in TREC 2010 Web Track. It can be seen from Figure 1(a) that this query has four official intents and intent  $i_1$  and  $i_3$  are related to the trade involving bobcat tractors. So we create a new node  $n_1$  that stands for "bobcat tractors" as their parent node. Similarly,  $n_1$  and  $i_4$  are related to "bobcat company," hence we create another new node  $n_2$  representing "bobcat company" as their parent. Finally, since  $n_2$  ("bobcat company") and  $i_2$  ("wild bobcat") are two distinct interpretations of query "bobcat," they are considered as the child nodes of the root node. The resultant OIH is shown in solid boxes in the left of Figure 1(b). Further, we extend the OIH by adding a child to  $i_2$ , and adding a child and a grandchild to  $i_4$ . The resultant EIH is shown in solid boxes plus dashed boxes in the left of Figure 1(b).

As for the OIH or EIH shown in Figure 1(b): (1) It is in a tree structure (Property 1); (2) Its root is query "bobcat" itself (Property 2); (3) The query is ambiguous, so the child nodes of root are its two different interpretations, i.e. "bobcat company" and "wild bobcat." A parent node refers to a more general concept than its children (Property 3), e.g. "bobcat company" is more general than "bobcat company homepage;" (4) The leaf nodes are exactly the official intents of query "bobcat" (Property 4). Further, the depth of all the leaf nodes in EIH is three (Property 5).

Only the relevance assessments for the original intents are available in TREC Web Track diversity test collections. For the intent hierarchies we create, document relevance judgments are just available for their leaf intents. As assessing document relevance is usually very time-consuming, it is not desirable to reassess the documents for the intermediate nodes. Fortunately, according to Property 3, a parent node of an intent hierarchy stands for a more general concept than its child nodes. Hence it is reasonable to assume that if a document is relevant to a node, it would be relevant to the node's parent. This means that we can



(a) Official intents of the query “defender”.



(b) LEFT: OIH is comprised of the solid boxes, whereas EIH includes both solid and dashed nodes. RIGHT: 1st column: document IDs ( $d$  : the ideal one), each of which is a ranked list of length 1. 2nd to 5th column (R and R') : relevance assessments for the official intents (in red) and derived relevance assessments for added nodes (in blue). 6th to 12th column (S): the scores computed by a measure at rank 1.

Fig. 2. The official intents, original intent hierarchy (OIH), and extended intent hierarchy (EIH) of No. 20 query “defender” in TREC Web Track 2009.

derive relevance assessments for the intermediate nodes starting from the leaves. In this paper, we simply let:

$$L_d(n) = \max_{c \in C(n)} L_d(c) \tag{9}$$

where  $L_d(n)$  is the relevance rating assigned to document  $d$  for node  $n$ , and  $C(n)$  is the set of child nodes of  $n$ .

We show an actual document (denoted by  $d$  in the following) from TREC Web Track 2010 diversity test collection in Figure 1(b). In the table, the officially provided relevance assessments are marked in blue, e.g. the relevance rating of  $d$  for  $i_1$  is 1. Firstly, node  $n_1$  has two child nodes,  $i_1$  and  $i_3$ , and the relevance ratings of  $d$  for them are 1 and 0. According to Equation (9), the relevance rating of  $d$  for  $n_1$  is 1. Similarly, we can derive the relevance rating for  $n_2$  based on its child node  $i_4$  and  $n_1$ . These derived relevance assessments are shown in red in the table of Figure 1(b).

To conclude, we create a new dataset containing intent hierarchies by manually grouping the official intents from TREC Web track test collections. The good news is that we do not need to reassess document relevance with regards to the intent hierarchies. We directly leverage document relevance assessments for the leaf intents, and automatically assign relevance ratings for the intermediate intents. This also implies that when we want to create hierarchical intents for evaluating diversity, we just need to assess document relevance for the leaf nodes or atomic intents.

Most of the time of creating the new dataset is spent on grouping the original intents. On average, we spend about three minutes per query mainly in understanding the original intents with the assistance of Google and Bing.

### 3.4 Hierarchical Measures

We assume that: (1) The user who is interested in an intent will be interested in the broader intents. Based on this

assumption, we automatically assign relevance assessment for the intermediate intents of an intent hierarchy according to Equation 9 in Section 3.3; (2) The user who is interested in an atomic intent will be likely to be more interested in the more related atomic intents, and be less interested in the less related atomic intents in an intent hierarchy. The relatedness between two atomic intents is defined by the length of the shortest path between the two intents: the shorter the length of the shortest path is, the more related the two intents are. This implies that the search results that cover more related atomic intents will only be of interested to a specific user group, and should be considered less diverse by the measures. N-rec and the measures based on N-rec reward search results that cover less related atomic intents in an intent hierarchy, which is illustrated in Section 3.4.1.

#### 3.4.1 Node Recall

Given a query  $q$ , let  $V$  denote the nodes in its intent hierarchy except for its root. Let  $d_r$  denote the document at rank  $r$ , and let  $N(d_r)$  denote the nodes in  $V$  to which  $d_r$  is relevant. Node recall ( $N-rec$ ) at rank  $K$  is defined as:

$$N-rec@K = \frac{\sum_{r=1}^K N(d_r)}{|V|} \tag{10}$$

N-rec is a natural generalization of I-rec when using the intent hierarchy. They both are rank-insensitive and cannot handle graded relevance assessments. I-rec credits minor intents, while N-rec credits minor nodes, which contribute to rewarding wide coverage of users’ information needs.

We use an example to show that N-rec outperforms I-rec in terms of discriminative power. In the right of Figure 2(b), I-rec<sub>1</sub>@1 means only using the first layer, I-rec<sub>2</sub>@1 means only using the second layer, and N-rec<sub>E</sub>@1 means using EIH when computing N-rec. These measures are computed at

rank 1. Note that the original I-rec is equal to I-rec<sub>2</sub>. We find that  $d_1 > d_2 = d_3$  according to I-rec<sub>1</sub>@1,  $d_1 = d_2 > d_3$  according to I-rec<sub>2</sub>@1, whereas  $d_1 > d_2 > d_3$  according to N-rec<sub>E</sub>@1. As we discussed in Section 3.4.2, The real preference is  $d_1 > d_2 > d_3$ . I-rec<sub>1</sub>@1 fails to tell the difference between  $d_2$  and  $d_3$ , while I-rec<sub>2</sub>@1 fails to distinguish between  $d_1$  and  $d_2$ . Only N-rec<sub>E</sub>@1 can tell the difference between the three documents, and thus is more discriminative than I-rec. The ranking list that covers less related atomic intents is rewarded by N-rec, which is consistent with the second user assumption.

Another point worth noting is that **the types of intent hierarchies are crucial to N-rec**. In the right of Figure 2(b), N-rec<sub>O</sub>@1 means using OIH instead of EIH. We find that N-rec<sub>O</sub>@1 cannot determine which one of  $d_1$  and  $d_2$  is better because they have exactly the same score. This indicates that using EIH has higher discriminative power than using OIH.

We aim to retrieve documents that cover as many nodes of intent hierarchies as possible. Besides, we prefer the documents that are highly relevant to more popular nodes and layers. N-rec mainly rewards wide coverage of different nodes of intent hierarchies in the top ranks. In the following, we will discuss some measures to complement N-rec.

### 3.4.2 Layer-Aware measures

Our next proposal is to first evaluate a ranked list for each layer of an intent hierarchy using existing measures, then combine the evaluation scores across multiple layers.

Let  $H$  denote the height of the intent hierarchy, and let  $L = l_1, l_2, \dots, l_H$  denote its  $H$  layers. We define Layer-Aware measures (*LA measures*) at document cutoff  $K$  as:

$$M-LA@K = \sum_{i=1}^H w_i M_i@K \quad (11)$$

Here,  $w_i$  is the weight of layer  $l_i$ , where  $\sum_{i=1}^H w_i = 1$ , and  $M_i$  is the evaluation score of measure  $M$  by using the intents of layer  $l_i$ . For example, ERR-IA-LA is computed as follows: (1) For each layer, compute the per-layer scores of ERR-IA; (2) Compute the weighted average of the per-layer scores using Equation (11).

We find that the combination of measures over layers of intent hierarchies could outperform the original measures. We use the query “defender”, No. 20 topic in TREC Web Track 2009 [29], as an example. We choose this query because it has a relatively simple intent hierarchy. Its EIH is shown in the left of Figure 2(b). Suppose we have three documents,  $d_1$ - $d_3$ , and each of them can be viewed as a ranked list containing only one document. Their relevance assessments for the EIH are displayed in blue in the right of Figure 2(b). To save space, the nodes that receive no relevant documents within the documents are not displayed. Assume that  $d$  is the first document within the ideal rank list and it is relevant to every node displayed. In the right of Figure 2(b), D<sub>#</sub>-nDCG<sub>1</sub>@1 is the evaluation score when only using the first layer of the EIH, D<sub>#</sub>-nDCG<sub>2</sub>@1 means only using the second layer, and D<sub>#</sub>-nDCG-LA<sub>E</sub>@1 is the average of D<sub>#</sub>-nDCG<sub>1</sub>@1 and D<sub>#</sub>-nDCG<sub>2</sub>@1. The original D<sub>#</sub>-nDCG is equal to D<sub>#</sub>-nDCG<sub>1</sub>. We use the measures to score  $d_1$  to  $d_3$ , i.e. evaluating at document cutoff 1.

We show the evaluation results in Figure 2(b):  $d_1 > d_2 = d_3$  under D<sub>#</sub>-nDCG<sub>1</sub>@1,  $d_1 = d_2 > d_3$  under D<sub>#</sub>-nDCG<sub>2</sub>@1,

whereas  $d_1 > d_2 > d_3$  under D<sub>#</sub>-nDCG-LA<sub>E</sub>@1. Here, “>” means the former document is preferred compared with the latter when evaluating them at rank 1, and “=” means neither is preferred. The real preference is  $d_1 > d_2 > d_3$ : (1)  $d_1$  is more diversified than  $d_2$  because  $d_1$  covers both “windows defender” and “defender arcade game online,” while  $d_2$  only covers the former; (2)  $d_2$  is more diversified than  $d_3$  because  $d_2$  covers both “windows defender homepage” and “windows defender reports” while  $d_3$  just covers the former. Here, only D<sub>#</sub>-nDCG-LA<sub>E</sub>@1 is consistent with the real preference. D<sub>#</sub>-nDCG<sub>1</sub>@1 fails to tell the difference between  $d_2$  and  $d_3$ , whereas D<sub>#</sub>-nDCG<sub>2</sub>@1 fails to tell the difference between  $d_1$  and  $d_2$ . This indicates that the combination over layers has higher potential to reflect real user satisfaction than the use of a flat list of intents.

### 3.4.3 HD-measures

The global gain of an intent hierarchy at rank  $r$  is given by:

$$GG_h(r) = \sum_{i=1}^H w_i GG_i(r) \quad (12)$$

where  $w_i$  is the weight of layer  $l_i$  and  $GG_i(r)$  is the global gain for layer  $l_i$  at rank  $r$ . Let  $CGG_h(r) = \sum_{k=1}^r GG_h(k)$ , which is the cumulative global gain for the intent hierarchy at rank  $r$ . Further, let  $GG_h(r)$  and  $CGG_h(r)$  denote the global gain and the cumulative global gain for the intent hierarchy at rank  $r$  in the ideal ranked list. The ideal list is obtained by listing up all the judged documents in descending order of global gains for the intent hierarchy. Let  $J(r) = 1$  if the document at rank  $r$  is relevant to the intent hierarchy, and  $J(r) = 0$  otherwise. Let  $C(r) = \sum_{k=1}^r J(k)$ . We define *HD-nDCG* and *HD-Q* at document cutoff  $K$  as:

$$HD-nDCG@K = \frac{\sum_{r=1}^K GG_h(r) = \log(r+1)}{\sum_{r=1}^K GG_h(r) = \log(r+1)} \quad (13)$$

$$HD-Q@K = \frac{1}{\min(K; R)} \sum_{r=1}^K J(r) \frac{C(r) + CGG_h(r)}{r + CGG_h(r)} \quad (14)$$

where  $R$  is the number of judged documents relevant to the intent hierarchy.

### 3.4.4 LD<sub>#</sub>-measures

We use the leaf nodes of intent hierarchies to compute D-measures, such as D-nDCG and D-Q. Then, *LD<sub>#</sub>-measure* is defined as:

$$LDJ\text{-measure}@K = N\text{-rec}@K + (1 - \gamma)D\text{-measure}@K \quad (15)$$

where  $\gamma$  is a parameter controlling the tradeoff between diversity and relevance. Since D-measures only use the leaf nodes, LD<sub>#</sub>-measures mainly reward high relevance with more popular leaves. Also, LD<sub>#</sub>-measures cannot handle the weights of layers. To tackle these, we propose HD<sub>#</sub>-measures and LAD<sub>#</sub>-measures in the next two sections.

### 3.4.5 HD<sub>#</sub>-measures

We define *HD<sub>#</sub>-measure* as:

$$HDJ\text{-measure}@K = N\text{-rec}@K + (1 - \gamma)HD\text{-measure}@K \quad (16)$$

where *HD-measure* can be HD-nDCG or HD-Q, and  $\gamma$  is a parameter between 0 and 1.

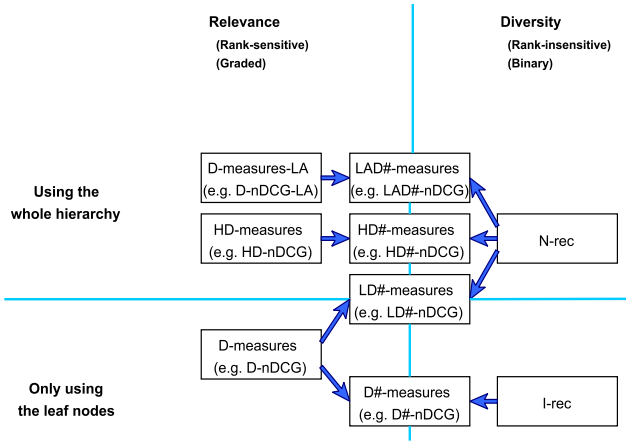


Fig. 3. Relationships of  $D_j$ -measures,  $LD_j$ -measures,  $HD_j$ -measures, and  $LAD_j$ -measures.

### 3.4.6 $LAD_j$ -measures

We define  $LAD_j$ -measure as:

$$LAD_j\text{-measure}@K = N\text{-rec}@K + (1 - \gamma)D\text{-measure-LA}@K \quad (17)$$

where  $\gamma$  is a parameter balancing diversity with relevance, and D-measure-LA is the LA version of D-measure.

To measure the relevance of ranked lists,  $HD_j$ -measures use HD-measures, while  $LAD_j$ -measures use D-measures-LA. HD-measures and D-measures-LA reward high relevance to more popular nodes, and can handle layer weights. The difference between them is what to combine over layers: HD-measures combine the global gain for each layer while D-measures-LA combine D-measures for each layer. Take HD-nDCG and D-nDCG-LA as an example:

$$HD\text{-nDCG}@K = \frac{\prod_{r=1}^K \left[ \prod_{i=1}^H w_i GG_i(r) \right]}{\prod_{r=1}^K \left[ \prod_{i=1}^H w_i GG_i(r) \right]}$$

$$D\text{-nDCG-LA}@K = \prod_{i=1}^H w_i D\text{-nDCG}_i@K$$

where  $GG_i(r)$  is the global gain for layer  $l_i$  at rank  $r$ , and  $D\text{-nDCG}_i$  means only using the nodes of layer  $l_i$ .

### 3.4.7 Summarization and Discussion

We call the proposed diversity measures *hierarchical measures*. LA measures mainly reward high relevance to major nodes, and thus take little account of minor nodes. This can be seen if we transform the definition of M-IA-LA, e.g., into

$$M\text{-IA-LA}@K = \prod_{n \in V} w_n M_n@K \quad (18)$$

where  $V$  is the nodes of an intent hierarchies except for the root, and  $w_n$  is the weight of node  $n$ .  $M_n$  is the per-node version of measure  $M$ , which can be nDCG, ERR, etc. This problem does not exist for the hierarchical measures which rely on N-rec because N-rec can credit minor intents.

Each of  $D_j$ -measures,  $LD_j$ -measures,  $HD_j$ -measures, and  $LAD_j$ -measures is a linear combination of two measures: one mainly rewards search result diversity, whereas another mainly rewards the relevance. We show their relationships in Figure 3. The figure shows that: (1) To reward the diversity,  $LD_j$ -measures,  $HD_j$ -measures, and  $LAD_j$ -measures

TABLE 1  
The types of intent hierarchies in TREC Web Track 2009-2013 and NTCIR-11 IMine

	(A) TREC (Uniform)		(B) NTCIR (EIH)	
	OIH	EIH	Uniform	Nonuniform
Top-down	OUT	EUT	EUT	ENT
Bottom-up	OUB	EUB	EUB	ENB

use the whole intent hierarchy, but  $D_j$ -measures only use the leaves; (2) To reward the relevance,  $HD_j$ -measures and  $LAD_j$ -measures use the whole intent hierarchy, but  $D_j$ -measures and  $LD_j$ -measures only use the leaves.

Hierarchical measures work more easily with EIH than with OIH because they require that for every layer, the node weights sum up to 1. EIH can guarantee this no matter which weighing scheme is used, while OIH cannot. So one drawback of using OIH is that when the node weights of one layer do not sum up to 1, renormalization is required. Besides, some hierarchical measures work in a layer wise way, like LA measures and HD-measures, but the document relevance assessments for some atomic intents may not be considered in the deep layers when using OIH. This means that hierarchical measures using OIH may neglect some atomic intents, and be counterintuitive sometimes.

## 4 EXPERIMENTS

First, we show that hierarchical measures using EIH weighted bottom-up, i.e. computing the node weights upwards given the leaf node weights, have advantages over existing measures in discriminative power and intuitiveness. Then, we present the experimental results of comparing the performance of hierarchical measures using different intent hierarchies, i.e. OIH or EIH weighted top-down or bottom-up uniformly or nonuniformly. Last, we show the benefits of using nonuniform weights, which are usually more costly than simple uniform weights.

### 4.1 Settings

We experiment with diversity measures on TREC Web Track 2009-2013 diversity test collections and NTCIR-11 IMine test collection. We build a new test collections from the TREC test collections, which contains intent hierarchies and the relevance assessments. The new dataset is publicly available on<sup>4</sup>. Available types of intent hierarchies in TREC Web Track 2009-2013 and NTCIR-11 IMine are summarized in Table 1. TREC Web Track 2009-2013 do not provide nonuniform intent weights, so only uniformly top-down and uniformly bottom-up weighting schemes can be applied. The intent hierarchies in NTCIR-11 IMine are EIH already. Since NTCIR-11 IMine provides nonuniform intent weights, the other two weighting schemes can also be applied, i.e. nonuniformly top-down and nonuniformly bottom-up. In the following, we use subscripts to denote the types of intent hierarchies.

Except for concordance test, we use document cutoff  $K = 20$  for all the measures. In concordance test, we use document cutoff  $K = 10$  instead because: (1) When doing case studies in Section 4.2.3, the details of top 20 documents

4. <http://www.playbigdata.com/dou/heval/>

TABLE 2

Discriminative power of existing measures and LA measures based on the paired bootstrap test at  $\alpha = 0.05$ .

(A) TREC Web Track 2009-2013 diversity test collections.			
Existing measure	disc. power	LA measure	disc. power
-nDCG	58.11%	-nDCG-LA	58.32%
ERR-IA	53.26%	ERR-IA-LA	53.79%
nDCG-IA	54.63%	nDCG-IA-LA	55.37%
Q-IA	47.47%	Q-IA-LA	48.95%
DJ-nDCG	57.05%	DJ-nDCG-LA	57.47%
DJ-Q	56.21%	DJ-Q-LA	56.32%
(B) NTCIR-11 IMine test collection.			
measure	disc. power	LA measure	disc. power
-nDCG	65.63%	-nDCG-LA	63.13%
ERR-IA	62.50%	ERR-IA-LA	65.00%
nDCG-IA	68.75%	nDCG-IA-LA	73.13%
Q-IA	65.00%	Q-IA-LA	71.88%
DJ-nDCG	68.75%	DJ-nDCG-LA	69.38%
DJ-Q	70.63%	DJ-Q-LA	71.25%

6

TABLE 3

Discriminative power of DJ-measures, LDJ-measures, HDJ-measures, and LADJ-measures based on the paired bootstrap test at  $\alpha = 0.05$ .

(A) TREC Web Track 2009-2013 diversity test collections.			
measure	disc. power	measure	disc. power
DJ-nDCG	57.05%	DJ-Q	56.21%
LDJ-nDCG	57.58%	LDJ-Q	56.42%
HDJ-nDCG	57.26%	HDJ-Q	56.53%
LADJ-nDCG	57.26%	LADJ-Q	56.53%
(B) NTCIR-11 IMine test collection.			
measure	disc. power	measure	disc. power
DJ-nDCG	68.75%	DJ-Q	70.63%
LDJ-nDCG	68.75%	LDJ-Q	71.25%
HDJ-nDCG	69.38%	HDJ-Q	71.25%
LADJ-nDCG	69.38%	LADJ-Q	71.25%

cannot be fitted into one page; (2) The conclusions drew when  $K = 10$  are almost the same as those drew when  $K = 20$ . We set  $\gamma = 0.5$  in Equation (7), (15), (16), and (17).

## 4.2 Hierarchical Measures Outperform Existing Measures

In this section, we fix the type of intent hierarchies used by hierarchical measures: (1) On TREC Web Track 2009-2013 test collections, hierarchical measures use EIH weighted uniformly bottom-up; (2) On NTCIR-11 IMine test collection, hierarchical measures use EIH weighted nonuniformly bottom-up. In the following sections, we will give explanation for selecting such types of intent hierarchies.

### 4.2.1 Discriminative Power

Following the previous work [17], [24], [26], [30], [31], we use the paired bootstrap test and set  $B = 1,000$  ( $B$  is the number of bootstrap samples). On TREC test collections we conduct the experiments as follows: (1) Sampling 20 submitted runs every year (2009-2013), which produces 950 pairs of sampled runs; (2) With the 950 pairs of sampled runs, computing the discriminative power using all the queries in TREC test collections. On the NTCIR test collection, we conduct the experiments as follows: (1) Using all the runs in NTCIR-11 IMine test collection because there are only 11 runs in Chinese and 15 runs in English. This gives us 160 pairs of runs; (2) With the 160 pairs of runs, computing the

TABLE 4

Intuitiveness based on preference agreement with gold standard measures. For each measure pair, the higher score is shown in bold.

(A) TREC Web Track 2009-2013. Gold standard: N-rec and Precision.				
	-nDCG-LA	ERR-IA-LA	nDCG-IA-LA	Q-IA-LA
LDJ-nDCG	<b>.673/.237</b>	.715/.185	<b>.385/.051</b>	<b>.375/.227</b>
HDJ-nDCG	<b>.667/.233</b>	.712/.181	<b>.392/.086</b>	<b>.372/.229</b>
LADJ-nDCG	<b>.672/.236</b>	<b>.714/.184</b>	<b>.385/.054</b>	<b>.374/.226</b>
LDJ-Q	<b>.787/.106</b>	<b>.814/.080</b>	<b>.678/.151</b>	<b>.578/.111</b>
HDJ-Q	<b>.787/.105</b>	<b>.814/.079</b>	<b>.678/.152</b>	<b>.578/.111</b>
LADJ-Q	<b>.787/.106</b>	<b>.814/.080</b>	<b>.678/.151</b>	<b>.578/.111</b>
(B) NTCIR-11 IMine. Gold standard: N-rec and Precision.				
	-nDCG-LA	ERR-IA-LA	nDCG-IA-LA	Q-IA-LA
LDJ-nDCG	<b>.705/.171</b>	<b>.806/.131</b>	<b>.653/.112</b>	<b>.654/.144</b>
HDJ-nDCG	<b>.704/.145</b>	<b>.813/.108</b>	<b>.624/.130</b>	<b>.633/.144</b>
LADJ-nDCG	<b>.721/.151</b>	<b>.815/.114</b>	<b>.657/.119</b>	<b>.657/.130</b>
LDJ-Q	<b>.785/.096</b>	<b>.835/.077</b>	<b>.718/.091</b>	<b>.728/.067</b>
HDJ-Q	<b>.792/.070</b>	<b>.845/.058</b>	<b>.718/.100</b>	<b>.736/.082</b>
LADJ-Q	<b>.782/.096</b>	<b>.837/.078</b>	<b>.720/.094</b>	<b>.730/.070</b>

discriminative power using all the queries in the NTCIR test collection. The results are shown in Tables 2 and 3.

By comparing the discriminative power of existing measures and the corresponding LA measures in each row of Table 2, we find that except  $\alpha$ -nDCG-LA, LA measures using EIH weighted bottom-up usually have higher discriminative power than (or have the same discriminative power as) the corresponding existing measures, especially in the case of IA measures. For example, (1) On TREC test collections, Q-IA-LA (49.16%) outperforms Q-IA (47.47%) in terms of discriminative power; (2) On the NTCIR test collection with nonuniform weights, Q-IA-LA (71.88%) beats Q-IA (65.00%) in terms of discriminative power.

By comparing the results of discriminative power of  $D_{\#}$ -measures,  $LD_{\#}$ -measures,  $HD_{\#}$ -measures, and  $LAD_{\#}$ -measures (each column in Table 3), we find that (1) When using EIH weighted bottom-up,  $LD_{\#}$ -measures,  $HD_{\#}$ -measures and  $LAD_{\#}$ -measures are consistently better than (or as good as)  $D_{\#}$ -measures in terms of discriminative power. For example, on TREC test collections,  $LD_{\#}$ -nDCG,  $HD_{\#}$ -nDCG, and  $LAD_{\#}$ -nDCG (57.58%, 57.26%, and 57.26% respectively) are all more discriminative than  $D_{\#}$ -nDCG (57.05%) in terms of discriminative power; (2) On the NTCIR test collection, the results of discriminative power of most measures are the same, which is due to the relatively small size of NTCIR-11 IMine test collection.

### 4.2.2 Concordance Test

In Section 3, we argue that LA measures may be less intuitive among hierarchical measures. On TREC test collections (or the NTCIR test collection), we compare the intuitiveness of two measures,  $M_1$  and  $M_2$ , as follows: (1) N-rec (or Precision) is used as the gold standard: The relative intuitiveness of  $M_1$  (or  $M_2$ ) is the ratio of ranked list pairs, for which  $M_1$  (or  $M_2$ ) and the gold standard have the same preference, to those for which  $M_1$  and  $M_2$  have different preference; (2) Both N-rec and Precision are used as the gold standard measures: The relative intuitiveness of  $M_1$  (or  $M_2$ ) is the ratio of ranked list pairs, for which  $M_1$  (or  $M_2$ ) has the same preference with N-rec and Precision, to those for which  $M_1$  and  $M_2$  have different preference. The results are shown in Table 4. N-rec is used as the gold standard measure for diversity because: (1) It is a simple binary measure; (2) It measures search result diversity better than I-rec, which is



TABLE 5

Intuitiveness based on preference agreement with gold standard measures. For each measure pair, the higher score is shown in bold.

(A) TREC Web Track 2009-2013. Gold standard: N-rec					
	ERR-IA	DJ-nDCG	LDJ-nDCG	HDJ-nDCG	LADJ-nDCG
-nDCG	<b>.987</b> /.390	.646/ <b>.979</b>	.632/ <b>.985</b>	.643/ <b>.980</b>	.640/ <b>.982</b>
ERR-IA	-	.561/ <b>.984</b>	.551/ <b>.990</b>	.559/ <b>.987</b>	.557/ <b>.988</b>
DJ-nDCG	-	-	.303/ <b>.741</b>	.646/ <b>.789</b>	.607/ <b>.794</b>
LDJ-nDCG	-	-	-	<b>.905</b> /.762	<b>.895</b> /.777
HDJ-nDCG	-	-	-	-	<b>.682</b> /.950
(B) TREC Web Track 2009-2013. Gold standard: Precision.					
	ERR-IA	DJ-nDCG	LDJ-nDCG	HDJ-nDCG	LADJ-nDCG
-nDCG	<b>.773</b> /.373	.493/ <b>.698</b>	.494/ <b>.698</b>	.492/ <b>.705</b>	.493/ <b>.704</b>
ERR-IA	-	.456/ <b>.725</b>	.456/ <b>.725</b>	.455/ <b>.731</b>	.455/ <b>.729</b>
DJ-nDCG	-	-	.580/ <b>.612</b>	.556/ <b>.692</b>	.560/ <b>.682</b>
LDJ-nDCG	-	-	-	.544/ <b>.736</b>	.547/ <b>.730</b>
HDJ-nDCG	-	-	-	-	<b>.762</b> /.529
(C) TREC Web Track 2009-2013. Gold standard: N-rec and Precision					
	ERR-IA	DJ-nDCG	LDJ-nDCG	HDJ-nDCG	LADJ-nDCG
-nDCG	<b>.761</b> /.082	.245/ <b>.681</b>	.236/ <b>.684</b>	.243/ <b>.686</b>	.241/ <b>.686</b>
ERR-IA	-	.187/ <b>.712</b>	.181/ <b>.715</b>	.186/ <b>.718</b>	.185/ <b>.717</b>
DJ-nDCG	-	-	.068/ <b>.372</b>	.300/ <b>.491</b>	.274/ <b>.486</b>
LDJ-nDCG	-	-	-	.457/ <b>.504</b>	.450/ <b>.513</b>
HDJ-nDCG	-	-	-	-	.456/ <b>.485</b>
(A) NTCIR-11 IMine. Gold standard: N-rec.					
	ERR-IA	DJ-nDCG	LDJ-nDCG	HDJ-nDCG	LADJ-nDCG
-nDCG	<b>.973</b> /.159	.385/ <b>.885</b>	.420/ <b>.867</b>	.458/ <b>.862</b>	.419/ <b>.883</b>
ERR-IA	-	.256/ <b>.944</b>	.284/ <b>.933</b>	.319/ <b>.922</b>	.292/ <b>.935</b>
DJ-nDCG	-	-	.548/ <b>.644</b>	<b>.643</b> /.622	.540/ <b>.698</b>
LDJ-nDCG	-	-	-	.719/ <b>.626</b>	.605/ <b>.721</b>
HDJ-nDCG	-	-	-	-	.475/ <b>.852</b>
(B) NTCIR-11 IMine. Gold standard: Precision.					
	ERR-IA	DJ-nDCG	LDJ-nDCG	HDJ-nDCG	LADJ-nDCG
-nDCG	<b>.886</b> /.227	.348/ <b>.841</b>	.466/ <b>.812</b>	.467/ <b>.829</b>	.467/ <b>.811</b>
ERR-IA	-	.287/ <b>.886</b>	.361/ <b>.868</b>	.363/ <b>.870</b>	.366/ <b>.860</b>
DJ-nDCG	-	-	<b>.788</b> /.587	<b>.720</b> /.636	<b>.762</b> /.619
LDJ-nDCG	-	-	-	.662/ <b>.727</b>	.663/ <b>.698</b>
HDJ-nDCG	-	-	-	-	<b>.738</b> /.639
(C) NTCIR-11 IMine. Gold standard: N-rec and Precision.					
	ERR-IA	DJ-nDCG	LDJ-nDCG	HDJ-nDCG	LADJ-nDCG
-nDCG	<b>.868</b> /.050	.100/ <b>.733</b>	.182/ <b>.698</b>	.213/ <b>.707</b>	.183/ <b>.710</b>
ERR-IA	-	.067/ <b>.834</b>	.116/ <b>.811</b>	.139/ <b>.802</b>	.120/ <b>.805</b>
DJ-nDCG	-	-	<b>.385</b> /.346	.413/.350	.365/.405
LDJ-nDCG	-	-	-	<b>.439</b> /.403	.337/ <b>.442</b>
HDJ-nDCG	-	-	-	-	.295/ <b>.525</b>

traditionally used as the gold standard measure for diversity. We find that LA measures are consistently less intuitive than LD $\ddagger$ -measures HD $\ddagger$ -measures, and LAD $\ddagger$ -measures.

We compare  $\alpha$ -nDCG, ERR-IA, D $\ddagger$ -nDCG, LD $\ddagger$ -nDCG, HD $\ddagger$ -nDCG, and LAD $\ddagger$ -nDCG in terms of intuitiveness. We do the concordance test on all the queries in TREC test collections or all the queries in the NTCIR test collection, and show the results in Table 5. In Table 5, blocks (A) and blocks (B) use N-rec and Precision as the gold standard measure respectively, whereas in blocks (C), both N-rec and Precision are used as the gold standard measures.

Table 5 shows that: (1) In terms of the diversity, LD $\ddagger$ -nDCG, HD $\ddagger$ -nDCG, and LAD $\ddagger$ -nDCG are usually more intuitive than existing measures. This is expected because these hierarchical measures depend on N-rec by means of Equation (15) and the like; (2) D $\ddagger$ -nDCG is consistently more intuitive than  $\alpha$ -nDCG and ERR-IA, but generally less intuitive than LD $\ddagger$ -nDCG, HD $\ddagger$ -nDCG, and LAD $\ddagger$ -nDCG; (3) Among the hierarchical measures, LD $\ddagger$ -nDCG is most intuitive in terms of diversity; HD $\ddagger$ -nDCG is most intuitive in terms of relevance; LAD $\ddagger$ -nDCG is the most intuitive

measure in terms of both diversity and relevance.

Table 5 shows that LD $\ddagger$ -nDCG, HD $\ddagger$ -nDCG and LAD $\ddagger$ -nDCG, which use the whole intent hierarchy to measure diversity, are more intuitive than D $\ddagger$ -nDCG in terms of diversity. HD $\ddagger$ -nDCG and LAD $\ddagger$ -nDCG, which use the whole intent hierarchy to measure relevance, are more intuitive than D $\ddagger$ -nDCG and LD $\ddagger$ -nDCG in terms of relevance. We get the same result when both diversity and relevance are considered. This indicates that **using the whole intent hierarchy instead of using the leaf nodes only can improve the intuitiveness of measures**. This is because: (1) Diversity is considered as covering as many nodes of an intent hierarchy as possible, so the measures using the whole intent hierarchy are more intuitive; (2) When using the whole intent hierarchy to compute diversity, using the whole intent hierarchy to compute relevance achieves a better balance between diversity and relevance than only using the leaves.

#### 4.2.3 Case Studies

D $\ddagger$ -nDCG, LD $\ddagger$ -nDCG, HD $\ddagger$ -nDCG, and LAD $\ddagger$ -nDCG are closely related (shown in Sections 3.4.7 and 4.2.4). Besides, they are consistently more intuitive than the others. We further examine their differences in terms of intuitiveness by looking at some real examples from the submitted runs in TREC Web Track 2009-2013 diversity task.

Specifically, we select five pairs of real ranked lists from TREC Web Track diversity runs in Table 6, and refer to them as **Case A-E**. For example, **Case A** stands for two runs cmuFuTop10D and THUIR10DvNov for No. 77 query; The middle column shows the relevance assessments of the top ten documents in each run (e.g. the first document retrieved by cmuFuTop10D is relevant to intent  $i_4$  with a relevance rating 1); The last four columns show the  $i$ 's for each query (e.g. score of cmuFuTop10D minus that of THUIR10DvNov) where arrows indicate which run has higher score under each measure. Note that in this section, the measures are computed for a document cutoff  $K = 10$  because we only have space to show top 10 documents in Table 6. We categorize five cases into two classes from the viewpoint of diversity (**Case A-C**) or relevance (**Case D-E**).

In **Case A**, we argue that D $\ddagger$ -nDCG is less intuitive than the other three. THUIR10DvNov covers both "bobcat company" and "wild bobcat" while cmuFuTop10D only covers the former (Please refer to the detailed description for the official intents of No. 77 query shown in Figure 1) although both runs cover three leaf intents. In this sense, THUIR10DvNov is more diversified than cmuFuTop10D and should be preferred. Note that this is also a case where I-rec cannot tell which run is better but N-rec can. The rightmost column of Table 6 shows that only D $\ddagger$ -nDCG disagrees with this intuition. In **Case B**, we argue that D $\ddagger$ -nDCG and HD $\ddagger$ -nDCG are less intuitive than the other two. Similar to **Case A**, UAMSD10aSRfu covers both "bobcat company" and "wild bobcat," whereas THUIR10DvQEW fails to cover the latter. So UAMSD10aSRfu should be preferred, and only LAD $\ddagger$ -nDCG and LD $\ddagger$ -nDCG agree with this. In **Case C**, we argue that LD $\ddagger$ -nDCG is the most intuitive among the four measures. In this case, both msrsv2div and qirdcsuog3 cover "bobcat company" and "wild bobcat". However, Figure 1 shows that msrsv2div covers both "bobcat tractors" and "bobcat company homepage," which are sub intents of

TABLE 6

Five ranked list pairs from TREC Web Track 2009-2013 diversity test collections. 1st column: case IDs (query IDs). 2nd column: run IDs. 3rd column: number of official intents covered by each run. 4th column: number of nodes in EIH covered by each run. 5th column: relevance ratings for each intent at ranks 1-10. The rightmost column: performance differences using each measure and arrows point to its preferred run.

				Document rank (i: official intents)										$\Delta$ in D $\ddagger$ - nDCG	$\Delta$ in LD $\ddagger$ - nDCG	$\Delta$ in HD $\ddagger$ - nDCG	$\Delta$ in LAD $\ddagger$ - nDCG
		1	2	3	4	5	6	7	8	9	10						
A (77)	cmuFuTop10D	3	6	$i_4 L1$	$i_3 L1$						$i_1 L1$	0.0013	-0.1098	-0.0977	-0.0988		
	THUIR10DvNov	3	8	$i_4 L1$	$i_1 L1$				$i_2 L1$				*	+	+	+	
B (77)	THUIR10DvQEW	2	5	$i_4 L1$										0.0300	-0.0256	0.0011	-0.0019
	UAMSD10aSRfu	2	6	$i_4 L1$				$i_2 L1$						*	+	*	+
C (77)	msrsv2div	3	8	$i_4 L1$	$i_2 L1$	$i_2 L1$	$i_2 L1$			$i_3 L1$	$i_2 L1$	-0.0329	0.0226	-0.0115	-0.0085		
	qirdcsuog3	3	7	$i_3 L1$	$i_1 L1$	$i_1 L1$			$i_1 L1$	$i_1 L1$	$i_2 L1$	+	*	+	+		
D (117)	qutir11a	3	5	$i_1 L1$	$i_2 L1$	$i_2 L1$	$i_1 L1$	$i_1 L2$	$i_2 L1$	$i_1 L2$	$i_1 L3$	$i_2 L1$	$i_2 L1$	-0.0030	-0.0030	0.0171	0.0148
	uwBBadhoc	3	5	$i_1 L3$	$i_3 L1$										+	+	*
E (128)	2011SiftR2	3	5	$i_1 L2$	$i_1 L2$								0.0087	0.0087	-0.0005	0.0004	
	UWatMDSdm	3	5	$i_1 L1$	$i_1 L1$	$i_1 L1$		$i_1 L1$		$i_1 L2$	$i_1 L2$		*	*	+	*	
				$i_2 L1$	$i_3 L1$												

TABLE 7

Kendall's  $\tau$  / Symmetric  $\tau_{ap}$  by averaging over TREC Web track 2009-2013 or NTCIR-11 IMine. Values ( $\geq .950$ ) are shown in bold.

(A) TREC Web Track 2009-2013 diversity test collections					
	ERR- IA	D $\ddagger$ - nDCG	LD $\ddagger$ - nDCG	HD $\ddagger$ - nDCG	LAD $\ddagger$ - nDCG
-nDCG	.923/.870	.840/.796	.845/.796	.843/.792	.844/.793
ERR-IA	-	.772/.699	.780/.706	.779/.704	.779/.706
D $\ddagger$ -nDCG	-	-	<b>.976/.959</b>	<b>.976/.957</b>	<b>.977/.960</b>
LD $\ddagger$ -nDCG	-	-	-	<b>.991/.988</b>	<b>.995/.993</b>
HD $\ddagger$ -nDCG	-	-	-	-	<b>.995/.994</b>
(B) NTCIR-11 IMine test collection					
	ERR- IA	D $\ddagger$ - nDCG	LD $\ddagger$ - nDCG	HD $\ddagger$ - nDCG	LAD $\ddagger$ - nDCG
-nDCG	.869/.763	.899/.922	.825/.843	.835/.878	.835/.878
ERR-IA	-	.787/.767	.713/.704	.723/.728	.723/.728
D $\ddagger$ -nDCG	-	-	.926/.913	.936/.949	.936/.949
LD $\ddagger$ -nDCG	-	-	-	<b>.990/.964</b>	<b>.990/.964</b>
HD $\ddagger$ -nDCG	-	-	-	-	<b>.999/.999</b>

“bobcat company,” while qirdcsuog3 does not cover “bobcat company homepage.” Because of this, msrsv2div should be preferred and only LD $\ddagger$ -nDCG agrees with this.

In summary, from the viewpoint of diversity, LD $\ddagger$ -nDCG is the most intuitive measure. HD $\ddagger$ -nDCG is less intuitive than LAD $\ddagger$ -nDCG, but is more intuitive than D $\ddagger$ -nDCG.

The two runs in Case D and in Case E have the same I-rec and N-rec, hence the measures’ preference is determined by their Precision part (e.g. D-nDCG if it is D $\ddagger$ -nDCG, and HD-nDCG if it is HD $\ddagger$ -nDCG). In Case D, we argue that D $\ddagger$ -nDCG and LD $\ddagger$ -nDCG are less intuitive than the other two. No matter whether measuring by I-rec or by N-rec, qutir11a and uwBBadhoc are equally good in terms of diversity. However, qutir11a should be preferred because its top ten documents are all relevant, whereas uwBBadhoc only has three. From the rightmost column of Table 6, we find that D $\ddagger$ -nDCG and LD $\ddagger$ -nDCG fail to reflect this. In Case E, we argue that HD $\ddagger$ -nDCG is the most intuitive. UWatMDSdm should be preferred because it returns much more relevant documents than 2011SiftR2. In this case, only HD $\ddagger$ -nDCG successfully recognizes this.

Generally, from the viewpoint of relevance, LAD $\ddagger$ -nDCG is more intuitive than LD $\ddagger$ -nDCG. LAD $\ddagger$ -nDCG is able to

measure the relevance of ranked lists more accurately by considering the whole intent hierarchy, and thus make the measures more consistent with Precision than LD $\ddagger$ -nDCG.

Case A serves as a good example showing that more diverse search results score higher under by hierarchical measures: THUIR10DvNov that covers both “bobcat company” and “wild bobcat” has higher score than cmuFuTop10D that only covers “bobcat company” under by LD $\ddagger$ -nDCG, HD $\ddagger$ -nDCG, as well as LAD $\ddagger$ -nDCG. Hierarchical measures achieve this by considering intents  $i_1$ ,  $i_3$ , and  $i_4$  are more related than intent  $i_2$  (see Figure 1 for the descriptions of the intents). Intuitively, hierarchical measures consider search result diversity as covering as many intents that are less related to each other as possible. This is the main difference between hierarchical measures and existing measures, which considers search result diversity as covering as many intents as possible

#### 4.2.4 Rank Correlation Results

We compute Kendall’s  $\tau$  and  $\tau_{ap}$  for different pairs of measures, and the results are shown in Table 7. We find that: (1) LD $\ddagger$ -nDCG, HD $\ddagger$ -nDCG, and LAD $\ddagger$ -nDCG are more correlated to D $\ddagger$ -nDCG than  $\alpha$ -nDCG and ERR-IA. This is because they are different kinds of extensions of D $\ddagger$ -nDCG. Similar to D $\ddagger$ -nDCG, they model diversity and relevance in different components separately. They yield the same evaluation results when the queries only have single-layer intent hierarchies; (3) LD $\ddagger$ -nDCG and HD $\ddagger$ -nDCG are less correlated. As discussed in Section 4.2.3, LD $\ddagger$ -nDCG prefers highly diversified ranked lists, whereas HD $\ddagger$ -nDCG prefers highly relevant ranked lists.

### 4.3 Using EIH is Better than Using OIH

The key factor influencing the performance of hierarchical measures is using OIH or using EIH. We have discussed some drawbacks of using OIH in Section 3.4. Since the intent hierarchies in NTCIR-11 IMine test collection have already satisfied the five properties of EIH, we only use TREC Web Track 2009-2013 test collections to compare the hierarchical measures using OIH and those using EIH.

TABLE 8

Discriminative power of diversity measures based on the paired bootstrap test at  $\alpha = 0.05$  on TREC Web Track 2009-2013 test collections.

OIH uniform top-down		EIH uniform top-down		OIH uniform bottom-up		EIH uniform bottom-up	
measure	disc.power	measure	disc.power	measure	disc.power	measure	disc.power
-nDCG-LA <sub>OUT</sub>	58.21%	-nDCG-LA <sub>EUT</sub>	<b>58.32%</b>	-nDCG-LA <sub>OUB</sub>	58.21%	-nDCG-LA <sub>EUB</sub>	<b>58.32%</b>
ERR-IA-LA <sub>OUT</sub>	52.32%	ERR-IA-LA <sub>EUT</sub>	<b>53.79%</b>	ERR-IA-LA <sub>OUB</sub>	52.32%	ERR-IA-LA <sub>EUB</sub>	<b>53.79%</b>
nDCG-IA-LA <sub>OUT</sub>	54.11%	nDCG-IA-LA <sub>EUT</sub>	<b>55.37%</b>	nDCG-IA-LA <sub>OUB</sub>	53.79%	nDCG-IA-LA <sub>EUB</sub>	55.16%
Q-IA-LA <sub>OUT</sub>	48.21%	Q-IA-LA <sub>EUT</sub>	48.95%	Q-IA-LA <sub>OUB</sub>	47.89%	Q-IA-LA <sub>EUB</sub>	<b>49.16%</b>
DJ-nDCG-LA <sub>OUT</sub>	56.32%	DJ-nDCG-LA <sub>EUT</sub>	<b>57.47%</b>	DJ-nDCG-LA <sub>OUB</sub>	55.68%	DJ-nDCG-LA <sub>EUB</sub>	57.37%
DJ-Q-LA <sub>OUT</sub>	54.42%	DJ-Q-LA <sub>EUT</sub>	56.32%	DJ-Q-LA <sub>OUB</sub>	54.53%	DJ-Q-LA <sub>EUB</sub>	<b>56.53%</b>
LDJ-nDCG <sub>OUT</sub>	57.16%	LDJ-nDCG <sub>EUT</sub>	<b>57.58%</b>	LDJ-nDCG <sub>OUB</sub>	56.63%	LDJ-nDCG <sub>EUB</sub>	57.37%
HDJ-nDCG <sub>OUT</sub>	56.53%	HDJ-nDCG <sub>EUT</sub>	<b>57.26%</b>	HDJ-nDCG <sub>OUB</sub>	56.21%	HDJ-nDCG <sub>EUB</sub>	<b>57.26%</b>
LADJ-nDCG <sub>OUT</sub>	56.53%	LADJ-nDCG <sub>EUT</sub>	<b>57.26%</b>	LADJ-nDCG <sub>OUB</sub>	56.00%	LADJ-nDCG <sub>EUB</sub>	<b>57.26%</b>
LDJ-Q <sub>OUT</sub>	55.79%	LDJ-Q <sub>EUT</sub>	56.42%	LDJ-Q <sub>OUB</sub>	55.68%	LDJ-Q <sub>EUB</sub>	<b>56.63%</b>
HDJ-Q <sub>OUT</sub>	55.79%	HDJ-Q <sub>EUT</sub>	<b>56.53%</b>	HDJ-Q <sub>OUB</sub>	55.47%	HDJ-Q <sub>EUB</sub>	<b>56.53%</b>
LADJ-Q <sub>OUT</sub>	54.42%	LADJ-Q <sub>EUT</sub>	<b>56.53%</b>	LADJ-Q <sub>OUB</sub>	54.42%	LADJ-Q <sub>EUB</sub>	<b>56.53%</b>

TABLE 9

Intuitiveness results on TREC Web Track 2009-2013 test collections. For each measure pair (one measure using OIH and the other using EIH), the higher score is shown in bold.

	Uniform top-down		Uniform bottom-up	
	OIH	EIH	OIH	EIH
(A) Gold standard measure: I-rec				
-nDCG-LA	1.000	1.000	1.000	1.000
ERR-IA-LA	.663	<b>.729</b>	.663	<b>.724</b>
nDCG-IA-LA	.669	<b>.724</b>	.624	<b>.748</b>
Q-IA-LA	.662	<b>.698</b>	.653	<b>.696</b>
DJ-nDCG-LA	.788	<b>.836</b>	.710	<b>.867</b>
DJ-Q-LA	.760	<b>.761</b>	.758	<b>.776</b>
LDJ-nDCG	.826	<b>.846</b>	.781	<b>.865</b>
HDJ-nDCG	.809	<b>.856</b>	.735	<b>.882</b>
LADJ-nDCG	.827	<b>.852</b>	.753	<b>.886</b>
LDJ-Q	.738	<b>.801</b>	.725	<b>.793</b>
HDJ-Q	.750	<b>.793</b>	.739	<b>.795</b>
LADJ-Q	<b>.814</b>	.761	<b>.815</b>	.774
(B) Gold standard measure: Precision				
-nDCG-LA	1.000	1.000	1.000	1.000
ERR-IA-LA	.539	<b>.655</b>	.539	<b>.655</b>
nDCG-IA-LA	.593	<b>.660</b>	.472	<b>.744</b>
Q-IA-LA	.610	<b>.659</b>	.510	<b>.743</b>
DJ-nDCG-LA	.562	<b>.659</b>	.520	<b>.666</b>
DJ-Q-LA	.509	<b>.718</b>	.496	<b>.730</b>
LDJ-nDCG	<b>.617</b>	.614	.557	<b>.672</b>
HDJ-nDCG	.577	<b>.655</b>	.517	<b>.668</b>
LADJ-nDCG	.569	<b>.655</b>	.511	<b>.676</b>
LDJ-Q	.631	<b>.633</b>	.602	<b>.631</b>
HDJ-Q	.620	<b>.649</b>	.584	<b>.648</b>
LADJ-Q	.500	<b>.761</b>	.483	<b>.754</b>
(C) Gold standard measure: I-rec and Precision				
-nDCG-LA	1.000	1.000	1.000	1.000
ERR-IA-LA	.349	<b>.481</b>	.349	<b>.481</b>
nDCG-IA-LA	.365	<b>.460</b>	.247	<b>.534</b>
Q-IA-LA	.379	<b>.451</b>	.304	<b>.510</b>
DJ-nDCG-LA	.390	<b>.510</b>	.298	<b>.545</b>
DJ-Q-LA	.298	<b>.486</b>	.285	<b>.510</b>
LDJ-nDCG	.462	<b>.473</b>	.366	<b>.545</b>
HDJ-nDCG	.406	<b>.520</b>	.296	<b>.556</b>
LADJ-nDCG	.415	<b>.526</b>	.307	<b>.568</b>
LDJ-Q	.374	<b>.436</b>	.333	<b>.426</b>
HDJ-Q	.375	<b>.444</b>	.330	<b>.445</b>
LADJ-Q	.323	<b>.503</b>	.309	<b>.529</b>

#### 4.3.1 Discriminative Power

On TREC test collections, we apply the same method as described in Section 4.2.1 to compute the discriminative power of hierarchical measures using four types of intent hierarchies, i.e. OUT, EUT, OUB, and EUB (shown in Table 1(A)). The results of discriminative power are shown in Table 8. By comparing the first column with the second column, and comparing the third column with the fourth column in Table 8, we find that the hierarchical measures using EIH are consistently more discriminative than those using OIH. This finding holds true no matter which weight-

ing scheme is used (UT or UB).

#### 4.3.2 Concordance Test

To compare the intuitiveness of hierarchical measures using OIH and those using EIH, we use I-rec and Precision as the gold standard measures. This is because: (1) They are simple binary measures, and are traditionally used as the gold standard measures; (2) They do not depend on intent hierarchies being OIH or EIH. We use all the queries in TREC test collections to compute the intuitiveness, and show the results in Table 9. By comparing the first column with the second column, and comparing the third column with the fourth column in Table 9, we find that the hierarchical measures using EIH are mostly more intuitive than those using OIH regardless of the weighting scheme used. This is because the hierarchical measures using OIH may reward high relevance to some official intents, and fail to reward wide coverage of the official intents.

To sum up, we find that hierarchical measures using EIH outperform those using OIH in terms of discriminative power and intuitiveness. This means that **hierarchical measures work better with EIH than with OIH**. This is because for an atomic intent, every layer of EIH either include it or one of its ancestors, but this is untrue for some layers of OIH. This makes hierarchical measures bias towards some atomic intents and be counterintuitive sometimes. For this reason, we use EIH when computing hierarchical measures on TREC test collections in Section 4.2.

#### 4.4 Weighting Bottom-up is Better than Weighting Top-down

Another factor that affects the performance of hierarchical measures is weighting intent hierarchies top-down or bottom-up. We conduct the experiments on TREC Web Track 2009-2013 test collections and NTCIR-11 IMine test collection, and find that the conclusions are consistent.

##### 4.4.1 Discriminative Power

We compute the discriminative power of hierarchical measures using different types of intent hierarchies on TREC Web Track test collections and NTCIR-11 IMine test collection, and show the results in Tables 8 and 10. We find that: (1) Hierarchical measures using OIH weighted top-down are mostly more discriminative than those using OIH weighted bottom-up; (2) Hierarchical measures using EIH weighted top-down are generally less discriminative than those using EIH weighted bottom-up.

TABLE 10  
Discriminative power of diversity measures based on the paired bootstrap test at  $\alpha = 0.05$  on NTCIR-11 IMine test collection.

EIH uniform top-down		EIH nonuniform top-down		EIH uniform bottom-up		EIH nonuniform bottom-up	
measure	disc.power	measure	disc.power	measure	disc.power	measure	disc.power
-nDCG-LA <sub>EUT</sub>	<b>63.13%</b>	-nDCG-LA <sub>ENT</sub>	<b>63.13%</b>	-nDCG-LA <sub>EUB</sub>	<b>63.13%</b>	-nDCG-LA <sub>ENB</sub>	<b>63.13%</b>
ERR-IA-LA <sub>EUT</sub>	<b>65.00%</b>	ERR-IA-LA <sub>ENT</sub>	<b>65.00%</b>	ERR-IA-LA <sub>EUB</sub>	<b>65.00%</b>	ERR-IA-LA <sub>ENB</sub>	<b>65.00%</b>
nDCG-IA-LA <sub>EUT</sub>	70.63%	nDCG-IA-LA <sub>ENT</sub>	71.25%	nDCG-IA-LA <sub>EUB</sub>	72.50%	nDCG-IA-LA <sub>ENB</sub>	<b>73.13%</b>
Q-IA-LA <sub>EUT</sub>	69.38%	Q-IA-LA <sub>ENT</sub>	68.75%	Q-IA-LA <sub>EUB</sub>	<b>71.88%</b>	Q-IA-LA <sub>ENB</sub>	<b>71.88%</b>
DJ-nDCG-LA <sub>EUT</sub>	68.75%	DJ-nDCG-LA <sub>ENT</sub>	69.38%	DJ-nDCG-LA <sub>EUB</sub>	<b>70.00%</b>	DJ-nDCG-LA <sub>ENB</sub>	69.38%
DJ-Q-LA <sub>EUT</sub>	70.63%	DJ-Q-LA <sub>ENT</sub>	70.63%	DJ-Q-LA <sub>EUB</sub>	<b>71.25%</b>	DJ-Q-LA <sub>ENB</sub>	<b>71.25%</b>
LDJ-nDCG <sub>EUT</sub>	68.75%	LDJ-nDCG <sub>ENT</sub>	68.75%	LDJ-nDCG <sub>EUB</sub>	<b>69.38%</b>	LDJ-nDCG <sub>ENB</sub>	68.75%
HDJ-nDCG <sub>EUT</sub>	<b>69.38%</b>	HDJ-nDCG <sub>ENT</sub>	68.75%	HDJ-nDCG <sub>EUB</sub>	<b>69.38%</b>	HDJ-nDCG <sub>ENB</sub>	<b>69.38%</b>
LADJ-nDCG <sub>EUT</sub>	<b>69.38%</b>	LADJ-nDCG <sub>ENT</sub>	68.75%	LADJ-nDCG <sub>EUB</sub>	<b>69.38%</b>	LADJ-nDCG <sub>ENB</sub>	<b>69.38%</b>
LDJ-Q <sub>EUT</sub>	70.63%	LDJ-Q <sub>ENT</sub>	70.63%	LDJ-Q <sub>EUB</sub>	<b>71.25%</b>	LDJ-Q <sub>ENB</sub>	<b>71.25%</b>
HDJ-Q <sub>EUT</sub>	<b>71.25%</b>	HDJ-Q <sub>ENT</sub>	<b>71.25%</b>	HDJ-Q <sub>EUB</sub>	<b>71.25%</b>	HDJ-Q <sub>ENB</sub>	<b>71.25%</b>
LADJ-Q <sub>EUT</sub>	<b>71.25%</b>	LADJ-Q <sub>ENT</sub>	<b>71.25%</b>	LADJ-Q <sub>EUB</sub>	<b>71.25%</b>	LADJ-Q <sub>ENB</sub>	<b>71.25%</b>

TABLE 11  
Discriminative power of diversity measures based on the paired bootstrap test at  $\alpha = 0.05$  on NTCIR-11 IMine test collection.

EIH uniform top-down		EIH nonuniform top-down		EIH uniform bottom-up		EIH nonuniform bottom-up	
measure	disc.power	measure	disc.power	measure	disc.power	measure	disc.power
D-nDCG <sub>EUT</sub>	77.50%	D-nDCG <sub>ENT</sub>	77.50%	D-nDCG <sub>EUB</sub>	77.50%	D-nDCG <sub>ENB</sub>	77.50%
HD-nDCG <sub>EUT</sub>	77.50%	HD-nDCG <sub>ENT</sub>	<b>78.13%</b>	HD-nDCG <sub>EUB</sub>	76.88%	HD-nDCG <sub>ENB</sub>	76.25%
D-nDCG-LA <sub>EUT</sub>	76.88%	D-nDCG-LA <sub>ENT</sub>	<b>78.13%</b>	D-nDCG-LA <sub>EUB</sub>	76.88%	D-nDCG-LA <sub>ENB</sub>	77.50%
D-Q <sub>EUT</sub>	81.88%	D-Q <sub>ENT</sub>	81.88%	D-Q <sub>EUB</sub>	81.88%	D-Q <sub>ENB</sub>	81.88%
HD-Q <sub>EUT</sub>	<b>81.88%</b>	HD-Q <sub>ENT</sub>	<b>81.88%</b>	HD-Q <sub>EUB</sub>	81.25%	HD-Q <sub>ENB</sub>	81.25%
D-Q-LA <sub>EUT</sub>	<b>81.88%</b>	D-Q-LA <sub>ENT</sub>	<b>81.88%</b>	D-Q-LA <sub>EUB</sub>	81.25%	D-Q-LA <sub>ENB</sub>	81.25%

#### 4.4.2 Concordance Test

To compare the intuitiveness of hierarchical measures using intent hierarchies weighted top-down or bottom-up, I-rec and Precision are used as the gold standard measures for the same reasons as mentioned in 4.3.2. We use all the queries in TREC test collections or all the queries the NTCIR test collection to compute the intuitiveness. The results are shown in Table 12. By comparing the first column with the second column in Table 12, we find that hierarchical measures using OIH weighted top-down are generally more intuitive than those using OIH weighted bottom-up. However, by comparing the third column with the fourth column, comparing the fifth column with the sixth column, and comparing the seventh column with the eighth column in Table 12, we find that hierarchical measures using EIH weighted top-down are generally less intuitive than those using EIH weighted bottom-up.

To summarize, it is preferable for OIH to be weighted top-down, whereas it is preferable for EIH to be weighted bottom-up. This is due to that when weighting EIH top-down, the leaves of a subtree that has many leaves will end up with minor weights, and may be neglected by hierarchical measures. Because of this, intent hierarchies are weighted bottom-up on NTCIR-11 IMine test collection in Section 4.2 for they are EIH. Sections 4.3 and 4.4 suggest that when nonuniform intent weights are not available, it is preferable for hierarchical measures to use EIH weighted bottom-up. Since TREC Web Track 2009-2013 do not provide nonuniform intent weights, hierarchical measures use EIH weighted bottom-up on the test collections in Section 4.2.

#### 4.5 Benefits of Using Nonuniform Weights

LD $\ddagger$ -measures, HD $\ddagger$ -measures, and LAD $\ddagger$ -measures use N-rec to model diversity. Since N-rec does not take weights into consideration, weighting schemes have no effect on its discriminative power. D-measures, HD-measures, and

D-measures-LA model relevance and consider weights. We compare their discriminative power on NTCIR-11 IMine test collections with different weighting schemes, and show the results in Table 11. By comparing the first column with the second column and comparing the third column and the fourth column in Table 11, we find that D-measures, HD-measures, and D-measures-LA using nonuniform weights tend to have higher discriminative power than those using uniform weights.

#### 4.6 Hierarchical Measures and Diversification Algorithms

Hu et al. [32] proposes two hierarchical diversification algorithms HxQuAD and HPM2 that use hierarchical intents. These two algorithms are demonstrated to outperform traditional diversification algorithms xQuAD, PM2, TxQuAD, and TPM2 that use a flat list of intents. Hu et al. [32] automatically generates hierarchical intents whose height is two. xQuAD, PM2, TxQuAD, and TPM2 can use the first-layer intents only, the second-layer intents only, or all the intents in the first-layer and second-layer, which results in 12 different algorithms. Given two measures  $M_1$  and  $M_2$ , the percentage of run pairs that are significantly different in the total run pairs is computed as follows: (1) The topic set in TREC Web Track 2009-2012 diversity test collections is randomly partitioned into five equal sized subsets. We fix the specific partition of the original topic set afterwards. For each of the five subsets, the remaining four subsets are used as the training data, and  $M_1$  is used to decide the best configuration of an algorithm in the training data, which gives us the final run for the spared subset. In this way, we produce a final run for each of the two hierarchical algorithms and 12 traditional algorithms; (2) We use  $M_2$  to evaluate the final runs produced before, one for each algorithm. Two-tailed paired t-test with significant level equal 0.05 is then used to test whether the run of a hierarchical algorithm

TABLE 12

Intuitiveness results on TREC Web Track 2009-2013 and NTCIR-11 IMine. For each measure pair, the higher score is shown in bold.

	TREC OIH uniform		TREC EIH uniform		NTCIR-11 EIH uniform		NTCIR-11 EIH nonuniform	
	Top-down	Bottom-up	Top-down	Bottom-up	Top-down	Bottom-up	Top-down	Bottom-up
(A) Gold standard measure: I-rec								
-nDCG-LA	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
ERR-IA-LA	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
nDCG-IA-LA	<b>0.775</b>	0.613	<b>0.727</b>	0.685	0.653	<b>0.721</b>	0.669	<b>0.673</b>
Q-IA-LA	<b>0.731</b>	0.651	<b>0.722</b>	0.667	0.654	<b>0.681</b>	<b>0.681</b>	0.620
DJ-nDCG-LA	<b>0.967</b>	0.641	0.814	<b>0.877</b>	0.636	<b>0.821</b>	0.753	<b>0.771</b>
DJ-Q-LA	<b>0.855</b>	0.825	0.727	<b>0.840</b>	0.818	<b>0.864</b>	0.864	<b>0.909</b>
LDJ-nDCG	<b>0.894</b>	0.847	0.821	<b>0.893</b>	0.565	<b>0.768</b>	0.736	<b>0.753</b>
HDJ-nDCG	<b>0.961</b>	0.689	0.840	<b>0.871</b>	0.705	<b>0.795</b>	<b>0.813</b>	0.731
LADJ-nDCG	<b>0.966</b>	0.669	0.829	<b>0.881</b>	0.691	<b>0.819</b>	0.779	<b>0.785</b>
LDJ-Q	0.764	<b>0.867</b>	0.786	<b>0.813</b>	<b>0.882</b>	0.647	<b>0.912</b>	0.794
HDJ-Q	0.790	<b>0.849</b>	0.768	<b>0.840</b>	0.837	0.837	0.800	<b>0.855</b>
LADJ-Q	<b>0.884</b>	0.824	0.775	<b>0.839</b>	0.837	0.837	0.843	0.843
(B) Gold standard measure: Precision								
-nDCG-LA	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
ERR-IA-LA	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
nDCG-IA-LA	<b>0.746</b>	0.483	0.618	<b>0.662</b>	0.620	<b>0.697</b>	0.618	<b>0.687</b>
Q-IA-LA	<b>0.780</b>	0.497	0.634	<b>0.663</b>	0.598	<b>0.657</b>	0.576	<b>0.629</b>
DJ-nDCG-LA	<b>0.644</b>	0.526	0.633	<b>0.638</b>	0.649	<b>0.709</b>	0.620	<b>0.753</b>
DJ-Q-LA	<b>0.805</b>	0.468	<b>0.740</b>	0.711	0.864	<b>0.886</b>	0.864	<b>0.886</b>
LDJ-nDCG	<b>0.739</b>	0.505	0.566	<b>0.661</b>	0.464	<b>0.761</b>	0.545	<b>0.781</b>
HDJ-nDCG	<b>0.667</b>	0.506	0.635	<b>0.647</b>	<b>0.705</b>	0.659	0.684	<b>0.720</b>
LADJ-nDCG	<b>0.652</b>	0.513	0.635	<b>0.647</b>	0.691	<b>0.698</b>	0.674	<b>0.727</b>
LDJ-Q	<b>0.760</b>	0.613	0.671	<b>0.731</b>	0.412	<b>0.765</b>	0.735	<b>0.824</b>
HDJ-Q	<b>0.807</b>	0.495	<b>0.739</b>	0.712	<b>0.857</b>	0.816	<b>0.836</b>	0.782
LADJ-Q	<b>0.790</b>	0.485	<b>0.733</b>	0.725	<b>0.857</b>	0.816	<b>0.863</b>	0.804
(C) Gold standard measure: I-rec and Precision								
-nDCG-LA	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
ERR-IA-LA	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
nDCG-IA-LA	<b>0.554</b>	0.251	<b>0.436</b>	0.429	0.407	<b>0.502</b>	0.429	<b>0.502</b>
Q-IA-LA	<b>0.578</b>	0.307	<b>0.460</b>	0.428	0.413	<b>0.492</b>	0.437	<b>0.445</b>
DJ-nDCG-LA	<b>0.613</b>	0.257	0.467	<b>0.519</b>	0.377	<b>0.543</b>	0.446	<b>0.566</b>
DJ-Q-LA	<b>0.660</b>	0.310	0.487	<b>0.551</b>	0.727	<b>0.750</b>	0.750	<b>0.795</b>
LDJ-nDCG	<b>0.638</b>	0.370	0.408	<b>0.559</b>	0.152	<b>0.536</b>	0.354	<b>0.567</b>
HDJ-nDCG	<b>0.629</b>	0.268	0.486	<b>0.521</b>	0.438	<b>0.472</b>	<b>0.518</b>	0.497
LADJ-nDCG	<b>0.618</b>	0.263	0.476	<b>0.530</b>	0.403	<b>0.523</b>	0.477	<b>0.552</b>
LDJ-Q	<b>0.524</b>	0.483	0.458	<b>0.543</b>	0.353	<b>0.471</b>	0.647	0.647
HDJ-Q	<b>0.597</b>	0.359	0.508	<b>0.552</b>	<b>0.694</b>	0.653	0.636	0.636
LADJ-Q	<b>0.673</b>	0.314	0.510	<b>0.564</b>	<b>0.694</b>	0.653	<b>0.706</b>	0.647

is significantly better than that of a traditional algorithm. The significance test is performed between a hierarchical algorithm and a traditional algorithm, which gives us 24 run pairs in total. Finally, we compute the percentage of run pairs, which have significant difference, in the 24 run pairs for the measures  $M_1$  and  $M_2$ . The results are shown in Table 13, where the measures in the leftmost column are used as training measures and the measures in the topmost row are used as validation measures. We find that: (1) No matter which measure is used to tune parameters, it is easier to show the significant improvements of hierarchical algorithms when using hierarchical measures rather than D-nDCG. This means that in some cases, existing measures cannot find the advantages of hierarchical algorithms but hierarchical measures can; (2)  $\alpha$ -nDCG remains a good choice to tune parameters of hierarchical algorithms.

### 5 CONCLUSIONS AND FUTURE WORK

In this paper, we argued that flat lists are not expressive enough to model the relationships between user intents. In view of this, we introduced intent hierarchies with four different weighting schemes. Then we proposed hierarchical measures that could work with intent hierarchies, and illustrated their advantages over existing measures. We experimented with a new test collection based on TREC Web Track 2009-2013 diversity test collections, and the NTCIR-11 IMine

TABLE 13

The percentage of run pairs that pass the two-tailed paired t-test with significant level equal 0.05 when comparing the runs by hierarchical algorithms and those by traditional algorithms. The leftmost measures are used to tune algorithm parameters in the training data, and the resulting final runs are evaluated by the topmost measures.

	- nDCG	ERR- IA	DJ/ nDCG	LDJ/ nDCG	HDJ/ nDCG	LADJ/ nDCG
-nDCG	66.67%	58.33%	58.33%	62.50%	62.50%	62.50%
ERR-IA	66.67%	66.67%	33.33%	45.83%	45.83%	45.83%
DJ-nDCG	29.17%	12.50%	70.83%	75.00%	75.00%	75.00%
LDJ-nDCG	50.00%	25.00%	62.50%	70.83%	70.83%	70.83%
HDJ-nDCG	33.33%	12.50%	70.83%	75.00%	70.83%	70.83%
LADJ-nDCG	33.33%	12.50%	70.83%	75.00%	75.00%	70.83%

test collection. Our main experimental findings are: (1) Hierarchical measures can be more discriminative than existing measures which use flat lists of intents; (2) LDJ-nDCG should be used when the diversity of search results is more valued than the relevance, whereas HDJ-nDCG should be used when the relevance is more important. LADJ-nDCG is a better choice when diversity and relevance are equally important; (3) The performance of hierarchical measures depends on the types of intent hierarchies. When nonuniform weights are unavailable, it is preferable for hierarchical measures to use EIH weighted bottom-up uniformly. When nonuniform weights are available, it is preferable for hierarchical measures to use EIH weighted bottom-up nonuniformly; (4) The gain of using hierarchical intents to

diversify search results may be vague to existing measures based on flat lists. It is important to evaluate hierarchical diversification algorithms using hierarchical measures.

## ACKNOWLEDGMENTS

Zhicheng Dou is the corresponding author of the paper. This work was supported by the National Natural Science Foundation of China under Grant No. 61502501 and No. 61402155, and the National Key Basic Research Program (973 Program) of China under Grant No. 2014CB340403.

## REFERENCES

- [1] Z. Dou, R. Song, and J.-R. Wen, "A large-scale evaluation and analysis of personalized search strategies," in *WWW*, 2007.
- [2] B. J. Jansen, A. Spink, and T. Saracevic, "Real life, real users, and real needs: a study and analysis of user queries on the web," *Information Processing & Management*, 2000.
- [3] C. Silverstein, H. Marais, M. Henzinger, and M. Moricz, "Analysis of a very large web search engine query log," *SIGIR Forum*, 1999.
- [4] C. L. Clarke, M. Kolla, and O. Vechtomova, "An effectiveness measure for ambiguous and underspecified queries," in *ICTIR*, 2009.
- [5] R. Agrawal, S. Gollapudi, A. Halverson, and S. Jeong, "Diversifying search results," in *WSDM*, 2009.
- [6] J. Carbonell and J. Goldstein, "The use of MMR, diversity-based reranking for reordering documents and producing summaries," in *SIGIR*, 1998.
- [7] H. Chen and D. R. Karger, "Less is more: probabilistic models for retrieving fewer relevant documents," in *SIGIR*, 2006.
- [8] Z. Dou, S. Hu, K. Chen, R. Song, and J.-R. Wen, "Multi-dimensional search result diversification," in *WSDM*, 2011.
- [9] F. Radlinski and S. Dumais, "Improving personalized web search using result diversification," in *SIGIR*, 2006.
- [10] R. L. Santos, C. Macdonald, and I. Ounis, "Intent-aware search result diversification," in *SIGIR*, 2011.
- [11] X. Zhu, A. B. Goldberg, J. Van Gael, and D. Andrzejewski, "Improving diversity in ranking using absorbing random walks." in *HLT-NAACL*, 2007.
- [12] R. L. Santos, C. Macdonald, and I. Ounis, "Exploiting query reformulations for web search result diversification," in *WWW*, 2010.
- [13] V. Dang and W. B. Croft, "Diversity by proportionality: an election-based approach to search result diversification," in *SIGIR*, 2012.
- [14] V. Dang and W. B. Croft, "Term level search result diversification," in *SIGIR*, 2013.
- [15] T. Sakai, N. Craswell, R. Song, S. Robertson, Z. Dou, and C.-Y. Lin, "Simple evaluation metrics for diversified search results," in *EVI*, 2010.
- [16] C. L. Clarke, M. Kolla, G. V. Cormack, O. Vechtomova, A. Ashkan, S. Büttcher, and I. MacKinnon, "Novelty and diversity in information retrieval evaluation," in *SIGIR*, 2008.
- [17] T. Sakai and R. Song, "Evaluating diversified search results using per-intent graded relevance," in *SIGIR*, 2011.
- [18] C. L. Clarke, N. Craswell, I. Soboroff, and G. V. Cormack, "Overview of the trec 2010 web track," in *TREC*, 2010.
- [19] X. Wang, Z. Dou, T. Sakai, and J.-R. Wen, "Evaluating search result diversity using intent hierarchies," in *SIGIR*, 2016.
- [20] C. X. Zhai, W. W. Cohen, and J. Lafferty, "Beyond independent relevance: methods and evaluation metrics for subtopic retrieval," in *SIGIR*, 2003.
- [21] K. Järvelin and J. Kekäläinen, "Ir evaluation methods for retrieving highly relevant documents," in *SIGIR*, 2000.
- [22] O. Chapelle, D. Metzler, Y. Zhang, and P. Grinspan, "Expected reciprocal rank for graded relevance," in *CIKM*, 2009.
- [23] C. L. Clarke, N. Craswell, I. Soboroff, and A. Ashkan, "A comparative analysis of cascade measures for novelty and diversity," in *WSDM*, 2011.
- [24] T. Sakai, "Evaluating evaluation metrics based on the bootstrap," in *SIGIR*, 2006.
- [25] B. A. Carterette, "Multiple testing in statistical analysis of systems-based information retrieval experiments," *TOIS*, 2012.
- [26] T. Sakai, "Evaluation with informational and navigational intents," in *WWW*, 2012.
- [27] M. G. Kendall, "A new measure of rank correlation," *Biometrika*, 1938.
- [28] E. Yilmaz, J. A. Aslam, and S. Robertson, "A new rank correlation coefficient for information retrieval," in *SIGIR*, 2008.
- [29] C. L. Clarke, N. Craswell, and I. Soboroff, "Overview of the trec 2009 web track," in *TREC*, 2009.
- [30] T. Sakai, "Bootstrap-based comparisons of ir metrics for finding one relevant document," in *AIRS*, 2006.
- [31] T. Sakai and S. Robertson, "Modelling a user population for designing information retrieval metrics," in *EVI*, 2008.
- [32] S. Hu, Z. Dou, X. Wang, T. Sakai, and J.-R. Wen, "Search result diversification based on hierarchical intents," in *CIKM*, 2015.



**Xiaojie Wang** is a Ph.D. candidate in the School of Computing and Information Systems at The University of Melbourne, Australia. He received his B.S degrees in Applied Mathematics and Computer Science from Renmin University of China in 2016. His research interests include information retrieval, data mining, and machine learning.



**Ji-Rong Wen** is a Professor and the Dean in the School of Information at Renmin University of China (RUC). He was awarded the prestigious National 1000 Plan Expert of China. He is the director of Big Data Research Centre and is leading Big Data Analytics and Intelligence Lab in RUC. His main research interest lies on big data management, information retrieval (especially web search), data mining and machine learning.



**Zhicheng Dou** is an Associate Professor in the School of Information at Renmin University of China. He received his B.S. and Ph.D. degrees in Computer Science from the Nankai University in 2003 and 2008, respectively. He worked at Microsoft Research as a researcher from July 2008 to September 2014. His research interests include natural language processing, information retrieval, and data mining.

**Tetsuya Sakai** is a professor at the Department of Computer Science and Engineering, Waseda University, Japan. He is also an Associate Dean at the IT Strategies Division of Waseda, and a visiting professor at National Institute of Informatics. He obtained a Ph.D from Waseda in 2000. His previous employment includes Toshiba and Microsoft Research Asia.

**Rui Zhang** is a Professor and Reader in the School of Computing and Information Systems at The University of Melbourne, Australia. His research interest is big data and information management in general, particularly in areas of intent tracking, recommendation systems, spatial and temporal data analytics, moving object management, indexing techniques, high performance computing and data streams.