

Overview of the NTCIR-12 MobileClick Task

Makoto P. Kato
Kyoto University
kato@dl.kuis.kyoto-u.ac.jp

Tetsuya Sakai
Waseda University
tetsuyasakai@acm.org

Takehiro Yamamoto
Kyoto University
tyamamot@dl.kuis.kyoto-u.ac.jp

Virgil Pavlu
Northeastern University
vip@ccs.neu.edu

Hajime Morita
Kyoto University
morita@nlp.ist.i.kyoto-u.ac.jp

Sumio Fujita
Yahoo Japan Corporation
sufujita@yahoo-corp.jp

ABSTRACT

This is an overview of the NTCIR-12 MobileClick-2 task (a sequel to 1CLICK in NTCIR-9 and NTCIR-10). In the MobileClick task, systems are expected to output a concise summary of information relevant to a given query and to provide immediate and direct information access for mobile users. We designed two types of MobileClick subtasks, namely, iUnit ranking and summarization subtasks, in which four research teams participated and submitted ? runs. We describe the subtasks, test collection, and evaluation methods and then report official results for NTCIR-12 MobileClick.

1. INTRODUCTION

Current web search engines usually return a ranked list of URLs in response to a query. After inputting a query and clicking on the search button, the user often has to visit several web pages and locate relevant parts within those pages. While these actions require significant effort and attention, especially for mobile users, they could be avoided if a system returned a concise summary of relevant information to the query [9].

The NTCIR-12 MobileClick task (and its predecessors, 1CLICK tasks organized in NTCIR-9 [10] and NTCIR-10 [1]) aims to directly return a summary of relevant information and immediately satisfy the user without requiring a lot of interaction with the device. Unlike the 1CLICK tasks, we expect the output to be a two-layered summary where the first layer contains the most important information and an outline of additional relevant information, while the second layer contains detailed information that can be accessed by clicking on links in the first layer. As shown in Figure 1, for query “NTCIR-11”, a MobileClick system presents general information about NTCIR-11 and a list of core tasks in the first layer. When the “MobileClick” link is clicked by the user, the system shows text in the second layer that explains the topic of that link.

Textual output of the MobileClick task is evaluated based on *information units (iUnits)* rather than document relevance. The performance of a submitted system is scored higher if it generates summaries including more important iUnits. In addition, we require systems to minimize the amount of text the user has to read or, equivalently, the time she has to spend in order to obtain relevant information. Although these evaluation principles were also taken into account in the 1CLICK tasks, they are extended to two-layered summaries where users can read a summary in multiple ways. We assume a user model that reads different parts of the summary by probabilistically clicking on links and compute an evaluation metric based on the importance of iUnits read as well as the time spent to obtain them.



Figure 1: An application of the MobileClick task. A concise two-layered summary can fit a small screen of the mobile device, and can satisfy diverse information needs.

MobileClick-2 attracted ? research teams from ? countries: ?, ?, and ?. Table ? provides a list of NTCIR-12 MobileClick participants with the number of iUnit ranking and summarization submissions. The total number of submissions was ?.

One of the biggest changes from the previous round of MobileClick was the evaluation system: we finished all the evaluation processes before releasing test data, and have returned evaluation results right after run submissions at our website ¹. This might enable participants to improve their systems based on returned results. In addition, the reproducibility was highly improved since there is no need to conduct additional assessments for new submissions. Another new trial in MobileClick-2 was *leader board*, by which participants can see evaluation results of the others. We expected more participants and higher performances by enhancing the visibility of state-of-the-arts performances achieved so far.

The remainder of this paper is structured as follows. Section 2 describes the details of the iUnit ranking and summarization subtasks. Section 3 introduces a test collection consisting of queries, iUnits, and a document collection. Section 4 describes our evaluation methodology. Section 5 reports on the official evaluation results for both subtasks. Finally, Section 6 concludes this paper.

2. SUBTASKS

MobileClick-2 comprises iUnit ranking and summarization subtasks. This section explains the two types of subtasks, and their input, output, and evaluation methodology.

2.1 iUnit Ranking Subtask

¹<http://www.mobileclick.org/>

The iUnit ranking subtask is a task where systems are expected to rank a set of information pieces (iUnits) based on their importance for a given query. This subtask was devised to enable *compartmentalized* evaluation, where we could separately evaluate the performance of estimating important iUnits and summarizing iUnits into a two-layered summary.

We provided a set of queries, a set of iUnits, and documents from which the iUnits were extracted. Note that the set of iUnits included irrelevant iUnits, which participants should rank below the other iUnits. We then asked participants to submit, for each query, a list of iUnits that are ordered by their estimated importance. More concretely, we accept a tab-delimited-values (TSV) file as an iUnit ranking run, where the first line must be a simple system description, and each of the other lines must represent a single iUnit. Therefore, a run file should look like the one shown below:

Listing 1: Example of an iUnit ranking run

```
This is an example run file
qid uid score
qid uid score
....
```

where “qid” is a query ID, “uid” is a iUnit ID, and “score” is estimated importance of the iUnit. In many ways, the iUnit ranking runs are similar with TREC ad-hoc runs in that they are essentially a ranked list of the objects retrieved. Note that we did not use “score” values for evaluation, and used the order of iUnits in run files.

2.2 iUnit Summarization Subtask

The iUnit summarization subtask is defined as follows: Given a query, a set of iUnits, and a set of intents, generate a structured textual output. In MobileClick, more precisely, the output must consist of two layers. The first layer is a list of iUnits and links to the second layer, while the second layer consists of lists of iUnits. Each link must be one of the provided intents and be associated with one of the iUnit lists in the second layer. Each list of iUnits in the first and second layers can include at most X characters so that it fits ordinary mobile screen size. The length of links is counted, while symbols and white spaces are excluded. In MobileClick-2, X is set to 420 for English and 280 for Japanese.

Each run must be a XML file that satisfies a DTD shown below:

Listing 2: DTD for an iUnit summarization run

```
<!ELEMENT results (sysdesc, result*)>
<!ELEMENT sysdesc (#PCDATA)>
<!ELEMENT result (first, second*)>
<!ELEMENT first (iunit | link)*>
<!ELEMENT second (iunit)*>
<!ELEMENT iunit EMPTY>
<!ELEMENT link EMPTY>
<!ATTLIST result qid NMTOKEN #REQUIRED>
<!ATTLIST iunit uid NMTOKEN #REQUIRED>
<!ATTLIST link iid NMTOKEN #REQUIRED>
<!ATTLIST second iid NMTOKEN #REQUIRED>
```

where

- The XML file includes a [results] element as the root element;
- The [results] element contains exactly one [sysdesc] element;
- The [results] element also contains [result] elements, each of

which corresponds a two-layered summary and has a [qid] attribute;

- A [result] element contains a [first] element and [second] elements;
- The [first] element contains [iunit] and [link] elements;
- A [second] element has an attribute [iid], and contains [iunit] elements.
- An [iunit] element has an attribute [uid] (iUnit ID); and
- A [link] element has an attribute [iid] (intent ID), which identifies a [second] element to be linked.

Note that the same [iunit] element may appear multiple times, e.g. an iUnit may appear in the [first] element and two [second] elements.

An XML file example that satisfies the DTD is shown below:

Listing 3: Example of an iUnit summarization run

```
<?xml version="1.0" encoding="UTF-8" ?>
<results>
  <sysdesc>
    Organizer Baseline
  </sysdesc>
  <result qid="MC-E-0001">
    <first>
      <iunit uid="MC-E-0001-U001" />
      <iunit uid="MC-E-0001-U003" />
      <link iid="MC-E-0001-I006" />
      <iunit uid="MC-E-0001-U004" />
      <link iid="MC-E-0001-I002" />
    </first>
    <second iid="MC-E-0001-I006">
      <iunit uid="MC-E-0001-U011" />
      <iunit uid="MC-E-0001-U019" />
    </second>
    <second iid="MC-E-0001-I002">
      <iunit uid="MC-E-0001-U029" />
      <iunit uid="MC-E-0001-U021" />
    </second>
  </result>
</results>
```

3. TEST COLLECTION

The NTCIR-12 MobileClick test collection includes queries, iUnits, intents, and a document collection. We describe the details of those components in the following subsections.

3.1 Queries

The NTCIR-12 MobileClick test collection includes 100 English and 100 Japanese queries (see Appendix A for the complete lists). Unlike the MobileClick-1 task, we selected more ambiguous/underspecified, or short queries like the 1CLICK tasks held in the past NTCIR. This is because we opt to focus on queries that are often utilized in mobile devices, and to tackle the problem of diverse intents in searchers.

We used a Wider Planet toolbar log from April to July 2014 for obtaining real-users’ queries, and translated them into English and Japanese. We selected frequent queries that belong to either *CELEBRITY*, *LOCAL*, and *DEFINITION* categories. Questions posted on Yahoo! Japan Chiebukuro² were used to generate QA queries. Those query categories were also employed in the 1CLICK tasks, since they are frequently used by mobile users [2].

²<http://chiebukuro.yahoo.co.jp/>

The definition of those categories is shown below (numbers in the brackets indicate the number of queries in the category):

CELEBRITY (20) names of celebrities such as artists, actors, politicians, and athletes.

LOCAL (20) landmarks and facilities (*e.g.* “tokyo sky tree”), or entities with geographical constraints (*e.g.* “banks Kyoto”).

DEFINITION (40) ambiguous terms that are often input to know their definition.

QA (20) natural language questions

3.2 Documents

To provide participants with a set of iUnits for each query, we downloaded 500 top-ranked documents that were returned by Bing search engine³ in response to each query, from which we extracted iUnits as explained in the next subsection. This crawling was conducted from May 29 to June 1, 2016. As we failed to access some of the documents, the number of downloaded documents per query is fewer than 500. The average number of documents for English queries is 418 and that for Japanese queries is 442.

NTCIR participants can obtain this document collection after their registration, and utilize them to estimate the importance of each iUnit and intent probability, *etc.*

3.3 iUnits

Like the ICLICK tasks held in the past NTCIR, we used iUnits as a unit of information in the MobileClick task. iUnits are defined as *relevant*, *atomic*, and *dependent* pieces of information, where

- *Relevant* means that an iUnit provides useful factual information to the user;
- *Atomic* means that an iUnit cannot be broken down into multiple iUnits without loss of the original semantics; and
- *Dependent* means that an iUnit can depend on other iUnits to be relevant.

Please refer to the ICLICK-2 overview paper for the details of the definition [1]. Although iUnits can depend on other iUnits to be relevant according to our definition, we excluded depending iUnits in this round for simplicity.

As this work requires careful assessment lasting for a long time and consideration on the three requirements of iUnits, we decided not to use crowd-sourcing mainly due to low controllability and high education cost. We hired assessors for extracting iUnits by hand and kept the quality of extracted iUnits by giving timely feedback on their results. Assessors were asked to extract as many iUnits as possible within an hour, by using *iUnit Extractor*⁴, a Firefox plugin we developed. The screenshot of the tool for iUnit extraction is shown in Figure 2. Each assessor worked on different sets of queries.

The total number of iUnits is 2,317 (23.8 iUnits per query) for English queries and 4,169 (41.7 iUnits per query) for Japanese queries. Examples of iUnits for English queries are shown in Table 1.

³<https://datamarket.azure.com/dataset/bing/search>

⁴<https://addons.mozilla.org/ja/firefox/addon/iunit-extractor/>

3.4 Intents

We introduce the notion of *intents* to the MobileClick-2 task, which have been utilized in the NTCIR INTENT and IMine tasks [3, 7, 12]. An intent can be defined as either a specific interpretation of an ambiguous query (“Mac OS” and “car brand” for “jaguar”), or an aspect of a faceted query (“windows 8” and “windows 10” for “windows”). In this round, intents were taken into account in evaluating the importance of iUnits, and were used as candidates of links to the second layer in the iUnit summarization subtask.

In the NTCIR INTENT and IMine tasks, the organizers clustered subtopics to form intents, while we constructed intents by clustering iUnits as follows:

- (1) Cluster iUnits by using a clustering interface,
- (2) Give each cluster a label representing iUnits included in the cluster, and
- (3) Let each label of a cluster represents an intent.

We hired assessors for the manual iUnit clustering, in which two iUnits were grouped together if

- (1) They are information about the same interpretation of an ambiguous query or the same aspect of a faceted query, and
- (2) They are likely to be interesting for the same user.

The criteria used in the label selection are listed below:

- (1) The label of a cluster should be descriptive enough for users to grasp the iUnits included in the cluster, and
- (2) The label of a cluster should be often used as a query or anchor text for the included iUnits.

These clustering and labeling tasks were conducted on *Clusty*⁵, a Web system for clustering. The screenshot of this system is shown in Figure 3.

As a result, we obtained 4.48 intents per query on average in the English subtasks, while we obtained 4.37 intents per query on average in the Japanese subtasks.

Subsequently, we let 10 crowd sourcing workers vote whether each intent is important or not. This voting was carried out to estimate the intent probability, which is the probability of intents of users who input a particular query, as was conducted in the NTCIR INTENT and IMine tasks. The assessors were asked to vote for multiple intents if they believed that they were interested in the intent when they had a chance to search by the query. We normalized the number of votes for each intent by the total number of votes for a query, and let $P(i|q)$ denote the normalized one for i , which we call intent probability of intent i of query q . More precisely, $P(i|q) = n_{i,q}/n_{.,q}$ where $n_{i,q}$ is the number of votes intent i received, and $n_{.,q}$ is the total number of votes for query q .

3.5 iUnit Importance

The importance of each iUnit was evaluated in terms of each intent, and *global importance* was derived from the per-intent importance and intent probability.

We asked two assessors to assess each iUnit in terms of each intent, and evaluate the importance at a five-point scale: 0 (unimportant), 1, 2 (somewhat important), 3, and 4 (highly important). The assessors were instructed to evaluate the importance by assuming

⁵<https://github.com/mpkato/clusty>

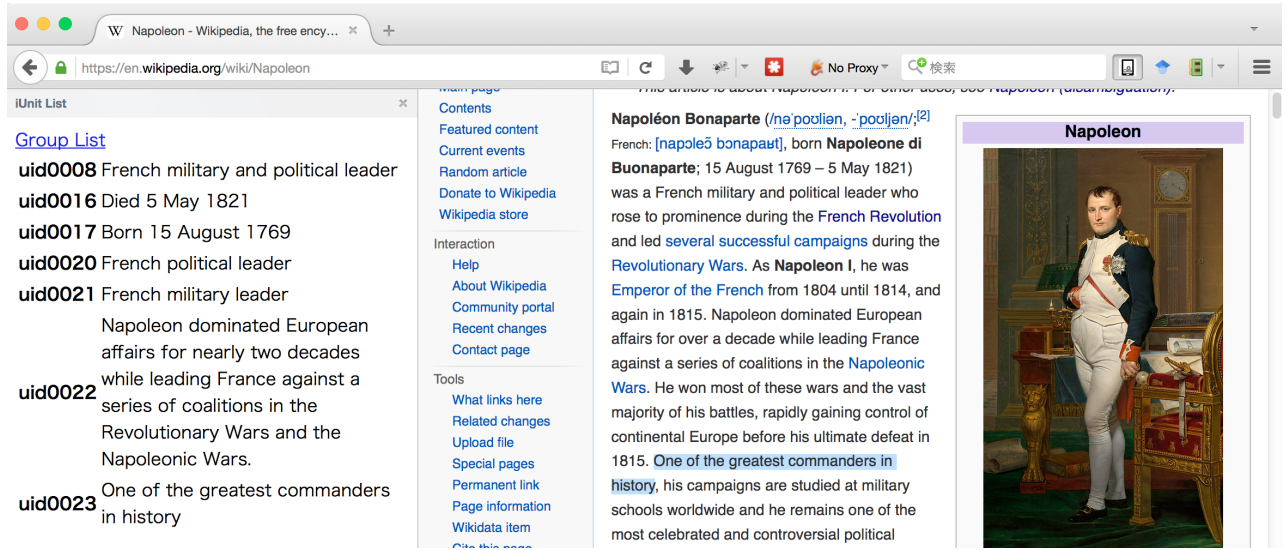


Figure 2: iUnit Extractor, a Firefox plugin for extracting iUnits from Web pages. Assessors can select a part of sentences and save its text, position, and URL by pressing Shift+Ctrl+x.

Table 1: Examples of iUnits for NTCIR-12 MobileClick English queries. Query MC2-E-0007 is “napoleon”.

Query ID	iUnit ID	iUnit
MC2-E-0007	MC2-E-0007-0001	born on the island of Corsica
MC2-E-0007	MC2-E-0007-0002	defeated at the Battle of Waterloo
MC2-E-0007	MC2-E-0007-0003	established legal equality and religious toleration
MC2-E-0007	MC2-E-0007-0004	an innovator
MC2-E-0007	MC2-E-0007-0005	absent during Peninsular War
MC2-E-0007	MC2-E-0007-0006	cut off European trade with Britain
MC2-E-0007	MC2-E-0007-0007	general of the Army of Italy
MC2-E-0007	MC2-E-0007-0008	one of the most controversial political figures
MC2-E-0007	MC2-E-0007-0009	won at the Battle of Wagram
MC2-E-0007	MC2-E-0007-0010	baptised as a Catholic

that they were interested in a given intent. We defined the importance of an iUnit in terms of an intent as follows: *an iUnit is more important if it is more necessary for more users who are interested in the intent*. For example, given intent “Mac OS” in response to query “jaguar”, iUnit “car company in UK” is *unimportant*, while it is *highly important* given intent “car brand”.

We used the average of the per-intent importance scores given by multiple assessors in our evaluation. The inter-assessor agreement was moderate: 0.556 in terms of *quadratic-weighted kappa* [11].

In the iUnit ranking subtask, we used the global importance of each iUnit for evaluation. Letting $P(i|q)$ be the intent probability of query q , the global importance of iUnit u is defined as follows:

$$G(u) = \sum_{i \in I_q} P(i|q) g_i(u), \quad (1)$$

where I_q is a set of intents for query q , and $g_i(u)$ denotes the per-intent importance of iUnit u in terms of intent i .

4. EVALUATION MEASURES

This section describes evaluation methodology used in the NTCIR-12 MobileClick tasks.

4.1 iUnit Ranking Subtask

Runs submitted by participants include a ranked list of iUnit IDs for each query, which can be handled in the same way as ad-hoc retrieval runs. Therefore, we employed standard evaluation metrics for ad-hoc retrieval in this subtask.

One of the evaluation metrics used in the iUnit ranking subtask was *normalized discounted cumulative gain* (nDCG). Discounted cumulative gain (DCG) is defined as follows:

$$\text{nDCG}@K = \sum_{r=1}^K \frac{G(u_r)}{\log_2(r+1)}, \quad (2)$$

where K is a cutoff parameter, and u_r is the r -th iUnit in a submitted ranked list. The normalized version of DCG (nDCG) is therefore defined as follows:

$$\text{nDCG}@K = \frac{\text{DCG}@K}{\text{iDCG}@K}, \quad (3)$$

where iDCG is DCG of the ideal ranked list of iUnits, which can be constructed by sorting all the iUnits for a query by their global importance.

MC2-E-0007

Download



Figure 3: Clusty, a Web system for clustering. Assessors can drag iUnits (blue rectangles) shown at the right pane, and drop them at one of the clusters shown at the left pane. In the screenshot, only “career” cluster was expanded to display all the iUnits in it.

Another evaluation metric is Q-measure proposed by Sakai [5]:

$$Q = \frac{1}{R} \sum_{r=1}^M \text{IsRel}(u_r) \frac{\sum_{r'=1}^r (\beta G(u_{r'}) + \text{IsRel}(u_{r'}))}{\beta \sum_{r'=1}^r G(u_{r'}^*) + r}, \quad (4)$$

where $\text{IsRel}(u)$ is an indicator function that returns 1 if $G(u) > 0$; otherwise 0, R is the number of iUnits with non-zero global importance (i.e. $\sum_u \text{IsRel}(u)$), M is the length of a ranked list, u_r^* is the r -th iUnit in the ideal ranked list of iUnits, and β is a patience parameter which we set to 1 following established standards [4]. Q-measure is used for ranking submitted runs since it can take into account the quality of the whole ranking.

Q-measure is a recall-based graded-relevance metric, while nDCG is a rank-based graded-relevance metric. Thus, we expect that using both metrics will enable us to measure the performance from different perspectives. Moreover, both of them were shown to be reliable [5].

4.2 iUnit Summarization Subtask

Runs submitted to the iUnit summarization subtask consists of the first layer \mathbf{f} and second layers $S = \{s_1, s_2, \dots, s_n\}$. The first layer \mathbf{f} consists of iUnits and links (e.g. $\mathbf{f} = (u_1, u_2, l_1, u_3)$ where u_j is an iUnit and l_j is a link). Each link l_j links to a second layer s_j . A second layer s_j is composed of iUnits (e.g. $s_1 = (u_{1,1}, u_{1,2}, u_{1,3})$).

The principles of the iUnit summarization evaluation metric are summarized as follows:

- (1) The evaluation metric is the expected utility of users who probabilistically read a summary.
- (2) Users are interested in one of the intents by following the intent probability $P(i|q)$.
- (3) Users read a summary following the rules below:

- (a) They read the summary from the beginning of the first layer in order and stop after reading L characters except symbols and white spaces.
- (b) When they reach the end of a link l_i , they click on the link and start to read its second layer if they are interested in the intent of l_i .
- (c) When they reach the end of a second layer s_j , they continue to read the first layer from the end of the link l_j .
- (4) The utility is measured by U-measure proposed by Sakai and Dou [6], which consists of a position-based gain and a position-based decay function.

We then generate the user tails (or *trailtext*) according to the user model explained above, compute a U-measure score for each trailtext, and finally estimate the expected U-measure by combining all the U-measure scores of different trailtexts. *M-measure*, an iUnit summarization evaluation metric, is defined as follows:

$$M = \sum_{\mathbf{t} \in T} P(\mathbf{t})U(\mathbf{t}), \quad (5)$$

where T is a set of all possible trailtexts, $P(\mathbf{t})$ is a probability of going through a trail \mathbf{t} , and $U(\mathbf{t})$ is the U-measure score of the trail. The computation of M-measure is illustrated in Figure 4.

A trailtext is a concatenation of all the texts read by a user, and can be defined as a list of iUnits and links in our case. According to our user model, a trailtext of a user who are interested in intent i can be obtained by inserting after the link of i a list of iUnits in its second layer. More specifically, trailtext \mathbf{t} of intent i is obtained as follows:

- (1) Let $\mathbf{f} = (\dots, u_{j-1}, l_k, u_j, \dots)$ where l_k is a link of intent i .
- (2) Generate $\mathbf{t} = (\dots, u_{j-1}, l_k, u_{k,1}, \dots, u_{k,|s_k|}, u_j, \dots)$ for second layer $s_k = (u_{k,1}, \dots, u_{k,|s_k|})$.

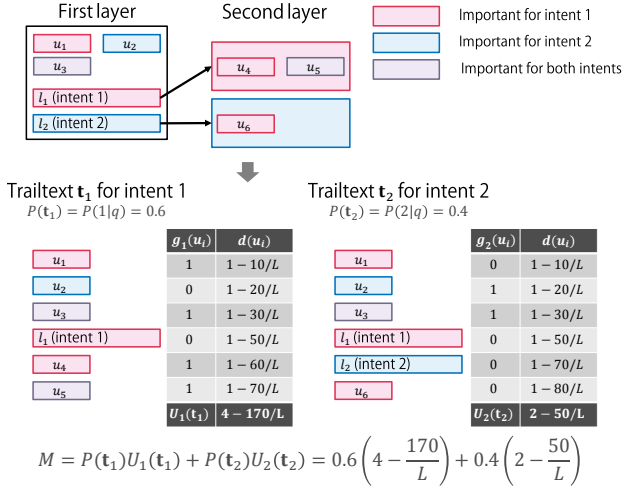


Figure 4: Illustration of M-measure computation. Two trailtexts are generated from a two-layered summary. U-measure of each trailtext is computed and summed with the probability of trailtexts (= intent probability). The length of iUnits is 10, while that of links is 20 in this example.

Note that a link in the trailtext is regarded as a non-relevant iUnit for the sake of convenience. Also note that only the first appearance of the same iUnit is relevant, while the other appearances are regarded as non-relevant.

As mentioned above, we can generate a trailtext for each intent, and do not need consider the other trailtexts as the way to read a summary only depends on the intent of users. In addition, the probability of a trailtext is equivalent to that of an intent for which the trailtext is generated. Thus, M-measure can be simply re-defined as follows:

$$M = \sum_{i \in I_q} P(i|q) U_i(t_i). \quad (6)$$

The U is now measured in terms of intent i in the equation above, since we assume that users going through t_i are interested in i .

The utility is measured by U-measure proposed by Sakai and Dou [6], and is computed by the importance and offset of iUnits in a trailtext. The offset of iUnit u in a trailtext is defined as the number of characters between the beginning of the trailtext and the end of u . More precisely, the offset of the j -th iUnit in trailtext t is defined as $\text{pos}_t(u) = \sum_{j'=1}^j \text{chars}(u_{j'})$ where $\text{chars}(u)$ is the number of characters of iUnit u except symbols and white spaces. Recall that a link in the trailtext contributes to the offset as a non-relevant iUnit. According to Sakai and Dou's work [6], U-measure is defined as follows:

$$U_i(t) = \frac{1}{\mathcal{N}} \sum_{j=1}^{|t|} g_i(u_j) d(u_j), \quad (7)$$

where d is a position-based decay function, and \mathcal{N} is a normalization factor (which we simply set to 1). The position-based decay function is defined as follows:

$$d(u) = \max \left(0, 1 - \frac{\text{pos}_t(u)}{L} \right), \quad (8)$$

where L is a patience parameter of users. Note that no gain can be obtained after L characters read, i.e. $d(u) = 0$. This is consistent with our user model in which users stop after reading L characters.

In MobileClick-2, L is set to twice as many as X : 840 for English and 560 for Japanese, since $L = 500$ (or 250) for Japanese was recommended by a study on S-measure [8].

5. RESULTS

Results will be reported after February 4, 2016.

6. CONCLUSIONS

This paper presents the overview of the MobileClick task at NTCIR-12. This task aims to develop a system that returns a concise summary of information relevant to a given query, and brings a structure into the summarization so that users can easily locate their desired information. In this paper, we mainly explained the task design, and evaluation methodology.

7. ACKNOWLEDGMENTS

We thank the NTCIR-12 MobileClick participants for their effort in submitting runs. We appreciate significant efforts made by Yahoo Japan Corporation for providing quite valuable search query data. We also thank Dr. Young-In Song from Wider Planet for providing useful data.

8. REFERENCES

- [1] M. P. Kato, M. Ekstrand-Abueg, V. Pavlu, T. Sakai, T. Yamamoto, and M. Iwata. Overview of the NTCIR-10 1CLICK-2 Task. In *NTCIR-10 Conference*, pages 243–249, 2013.
- [2] J. Li, S. Huffman, and A. Tokuda. Good abandonment in mobile and pc internet search. In *Proc. of SIGIR 2009*, pages 43–50, 2009.
- [3] Y. Liu, R. Song, M. Zhang, Z. Dou, T. Yamamoto, M. Kato, H. Ohshima, and K. Zhou. Overview of the NTCIR-11 IMine task. In *Proc. of NTCIR-11 Conference*, pages 8–23, 2014.
- [4] T. Sakai. On penalising late arrival of relevant documents in information retrieval evaluation with graded relevance. In *Proceedings of the First Workshop on Evaluating Information Access (EVIA 2007)*, pages 32–43, 2007.
- [5] T. Sakai. On the reliability of information retrieval metrics based on graded relevance. *Information processing & management*, 43(2):531–548, 2007.
- [6] T. Sakai and Z. Dou. Summaries, ranked retrieval and sessions: a unified framework for information access evaluation. In *Proc. of SIGIR 2013*, pages 473–482, 2013.
- [7] T. Sakai, Z. Dou, T. Yamamoto, Y. Liu, M. Zhang, R. Song, M. Kato, and M. Iwata. Overview of the ntcir-10 intent-2 task. *Proceedings of NTCIR-10*, pages 94–123, 2013.
- [8] T. Sakai and M. P. Kato. One click one revisited: Enhancing evaluation based on information units. In *Proc. of AIRS 2012*, pages 39–51, 2012.
- [9] T. Sakai, M. P. Kato, and Y.-I. Song. Click the search button and be happy: Evaluating direct and immediate information access. In *Proc. of CIKM 2011*, pages 621–630, 2011.
- [10] T. Sakai, M. P. Kato, and Y.-I. Song. Overview of NTCIR-9 1CLICK. In *Proceedings of NTCIR-9*, pages 180–201, 2011.
- [11] J. Sim and C. C. Wright. The kappa statistic in reliability studies: use, interpretation, and sample size requirements. *Physical therapy*, 85(3):257–268, 2005.
- [12] 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. *Proceedings of NTCIR-9*, pages 82–105, 2011.

APPENDIX

A. QUERIES

Full lists of English and Japanese queries used in MobileClick-2 are shown in Tables 2 and 3, respectively.

Table 2: NTCIR-12 MobileClick English queries.

ID	Query
MC2-E-0001	hulk hogan
MC2-E-0002	bruno mars
MC2-E-0003	ryan gigs
MC2-E-0004	sharon stone
MC2-E-0005	christopher nolan
MC2-E-0006	john jones
MC2-E-0007	napoleon
MC2-E-0008	selena gomez
MC2-E-0009	hemingway
MC2-E-0010	gareth bale
MC2-E-0011	tarantino
MC2-E-0012	darvish
MC2-E-0013	george clooney
MC2-E-0014	emma stone
MC2-E-0015	ronaldo
MC2-E-0016	park ji-sung
MC2-E-0017	robert downey jr
MC2-E-0018	millet
MC2-E-0019	miranda kerr
MC2-E-0020	emilia clarke
MC2-E-0021	bank adelanto
MC2-E-0022	cafe killeen
MC2-E-0023	cincinnati art museum
MC2-E-0024	delaware memorial bridge
MC2-E-0025	dermatology pittsburgh
MC2-E-0026	florida state university
MC2-E-0027	hotel oak creek
MC2-E-0028	junior high school sacramento
MC2-E-0029	kfc snohomish
MC2-E-0030	library cleveland tn
MC2-E-0031	mcdonald toole
MC2-E-0032	museum allentown
MC2-E-0033	nahanni national park
MC2-E-0034	orlando international airport
MC2-E-0035	orthopedic saint louis
MC2-E-0036	pharmacy nevada city
MC2-E-0037	pizza hut chichester
MC2-E-0038	post office olathe
MC2-E-0039	scottsdale stadium
MC2-E-0040	starbucks visalia
MC2-E-0041	bitcoin
MC2-E-0042	jaguar
MC2-E-0043	divers disease
MC2-E-0044	kebab
MC2-E-0045	windows 7
MC2-E-0046	april fool
MC2-E-0047	asiana airlines
MC2-E-0048	euro
MC2-E-0049	tetris
MC2-E-0050	nlb
MC2-E-0051	health insurance
MC2-E-0052	leica
MC2-E-0053	subarachnoid hemorrhage
MC2-E-0054	golf
MC2-E-0055	gold lacquer
MC2-E-0056	mango
MC2-E-0057	kakaotalk
MC2-E-0058	the age of discovery
MC2-E-0059	minecraft
MC2-E-0060	putty
MC2-E-0061	ups
MC2-E-0062	godzilla
MC2-E-0063	playstation
MC2-E-0064	smartphone
MC2-E-0065	aladdin
MC2-E-0066	pomeranian
MC2-E-0067	enka
MC2-E-0068	premier league
MC2-E-0069	quadcopter
MC2-E-0070	lettuce
MC2-E-0071	yellow ribbon
MC2-E-0072	real madrid
MC2-E-0073	baguette
MC2-E-0074	titanic
MC2-E-0075	accent
MC2-E-0076	yahoo!
MC2-E-0077	acai
MC2-E-0078	melon
MC2-E-0079	fuji xerox
MC2-E-0080	lineage
MC2-E-0081	difference between red and white wine
MC2-E-0082	aftershock mechanism
MC2-E-0083	how to cook coleslaw
MC2-E-0084	role of animal tail
MC2-E-0085	difference across to bcc cc
MC2-E-0086	difference across republic democratic republic federal
MC2-E-0087	way to reduce belly fat
MC2-E-0088	pc remote control method
MC2-E-0089	what is mirror made of
MC2-E-0090	why is white chocolate white
MC2-E-0091	why do children in costumes ask for candy at halloween
MC2-E-0092	difference in pepper types
MC2-E-0093	reason why dove is a symbol of peace
MC2-E-0094	etymology of japan
MC2-E-0095	how to make macaroons
MC2-E-0096	meaning of h and b pencil
MC2-E-0097	why is salt added when boiling spaghetti
MC2-E-0098	how to write a report
MC2-E-0099	principle of magic mirror
MC2-E-0100	difference of gdp and gnp

Table 3: NTCIR-12 MobileClick Japanese queries.

ID	Query
MC2-J-0001	ハルクホーガン
MC2-J-0002	八木アリサ
MC2-J-0003	ライアン・ギグス
MC2-J-0004	シャロン・ストーン
MC2-J-0005	ミーシャ
MC2-J-0006	ジョン・ジョーンズ
MC2-J-0007	ナポレオン
MC2-J-0008	セレナ・ゴメス
MC2-J-0009	ヘミングウェイ
MC2-J-0010	能年玲奈
MC2-J-0011	タランティノー
MC2-J-0012	ダルビッシュ
MC2-J-0013	東野圭吾
MC2-J-0014	エマ・ストーン
MC2-J-0015	ロナウド
MC2-J-0016	バク・チソン
MC2-J-0017	Yui
MC2-J-0018	ミレー
MC2-J-0019	ミランダ・カー
MC2-J-0020	エミリア・クラーク
MC2-J-0021	図書館明石市
MC2-J-0022	博物館豊橋市
MC2-J-0023	郵便局江ノ島
MC2-J-0024	とんかつ静岡市駅前
MC2-J-0025	kfc 宇都宮駅
MC2-J-0026	国立近代美術館
MC2-J-0027	マクドナルド博多駅
MC2-J-0028	最高裁判所
MC2-J-0029	ダイソー大宮駅周辺
MC2-J-0030	スターバックス岡山駅
MC2-J-0031	中学校東淀川区
MC2-J-0032	ロッテリア鳥取
MC2-J-0033	整形外科恵比寿
MC2-J-0034	ホテル勝浦駅
MC2-J-0035	カフェ川之江駅
MC2-J-0036	国会図書館
MC2-J-0037	銀行いすみ市
MC2-J-0038	皮膚科京都市北区
MC2-J-0039	高速バスターミナル札幌
MC2-J-0040	ビザハット川口市
MC2-J-0041	ビットコイン
MC2-J-0042	ジャガー
MC2-J-0043	潜水病
MC2-J-0044	ケバブ
MC2-J-0045	Windows 7
MC2-J-0046	エイプリル・フル
MC2-J-0047	アジアナ航空
MC2-J-0048	ユーロ
MC2-J-0049	テトリス
MC2-J-0050	nlb
MC2-J-0051	健康保険
MC2-J-0052	ライカ
MC2-J-0053	くも膜下出血
MC2-J-0054	ゴルフ
MC2-J-0055	蒔絵
MC2-J-0056	マンゴー
MC2-J-0057	カカオトーク
MC2-J-0058	大航海時代
MC2-J-0059	マインクラフト
MC2-J-0060	putty
MC2-J-0061	ups
MC2-J-0062	ゴジラ
MC2-J-0063	プレイステーション
MC2-J-0064	スマートフォン
MC2-J-0065	アラジン
MC2-J-0066	ボメラニアン
MC2-J-0067	演歌
MC2-J-0068	プレミアリーグ
MC2-J-0069	神の贈り物
MC2-J-0070	レタス
MC2-J-0071	黄色いリボン
MC2-J-0072	リアル・マドリッド
MC2-J-0073	バゲット
MC2-J-0074	タイタニック
MC2-J-0075	アクセント
MC2-J-0076	ヤフー
MC2-J-0077	アサイ
MC2-J-0078	メロン
MC2-J-0079	富士ゼロックス
MC2-J-0080	リネージュ
MC2-J-0081	赤ワインと白ワインの違い
MC2-J-0082	余震が起きるメカニズム
MC2-J-0083	コールスローを作る方法
MC2-J-0084	動物の尻尾の役割
MC2-J-0085	To、Bcc、Ccの違い
MC2-J-0086	共和国と民主共和国と連邦の違い
MC2-J-0087	おなかの脂肪を減らす方法
MC2-J-0088	コンピュータのリモート制御方法
MC2-J-0089	鏡は何からできているのか
MC2-J-0090	ホワイトチョコが白い理由
MC2-J-0091	ハロウィンでなぜ子供が変装してお菓子をもらうのか
MC2-J-0092	ペッパーの種類の違い
MC2-J-0093	鳩が平和の象徴とされる由来
MC2-J-0094	Japanの語源は何ですか
MC2-J-0095	マカロンの作り方
MC2-J-0096	鉛筆のHやBの意味
MC2-J-0097	スパゲティをゆでる時に塩を入れるのは何故
MC2-J-0098	レポートを書く方法
MC2-J-0099	マジックミラーの原理
MC2-J-0100	GDPとGNPの違い