

# 検索エンジンAPIを用いて収集したウェブページにおけるトピック多様性の分析

Analysis on Diversity of Topics in Web Pages collected using a Search Engine API

小池 大地\*<sup>1</sup> 宇津呂 武仁\*<sup>2</sup> 河田容英\*<sup>3</sup> 吉岡 真治\*<sup>4</sup>  
Daichi Koike Takehito Utsuro Yasuhide Kawada Masaharu Yoshioka

\*<sup>1</sup>筑波大学大学院システム情報工学研究科 Grad. Sch. Sys. & Inf. Eng, Univ of Tsukuba \*<sup>2</sup>筑波大学システム情報系 Fclty. Eng, Inf. & Sys, Univ of Tsukuba \*<sup>3</sup>(株) ログワークス Logworks Co., Ltd.

\*<sup>4</sup>北海道大学大学院情報科学研究科 Grad. Sch of Inf. Sci. & Tech, Hokkaido Univ.

This paper proposes how to collect Web pages of diverse topics using a search engine API, where a topic model is employed to estimate distribution of topics. In the proposed framework, given a certain query, relevant news articles and blog posts are collected and the topic distribution within the collected document set is estimated by a topic model. Next, from each topic within the document set consisting of news and blogs, a few query terms for the search engine API are automatically selected and are utilized in the procedure of collecting Web pages closely relevant to topics observed only in news or blogs. Experimental evaluation results show that the proposed procedure of collecting Web pages through queries collected from news and blogs achieve much more diverse distribution of topics compared with the baseline procedure of collecting Web pages through the search engine API without considering topics in news and blogs.

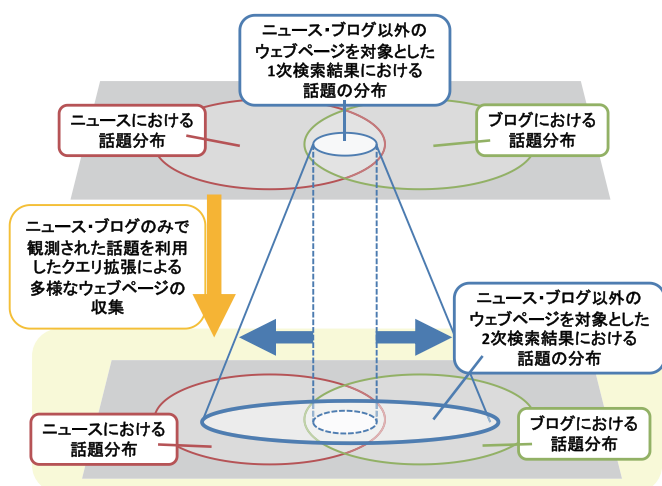


図 1: ニュース・ブログのみで観測された話題を利用したクエリ拡張による多様なウェブページの収集

## 1. はじめに

本論文では、検索エンジンAPIを用いたウェブページ収集タスクにおいて、トピックモデルを用いたトピックの分布の観点から、できる限り多様なウェブページを収集する方式を提案する。提案方式においては、まず、特定のクエリに関連するニュース・ブログを収集した文書集合を対象としてトピックモデルを適用し、トピックの分布を求める。そして、同一のクエリを対象としてニュース・ブログ以外のウェブページから収集した文書集合（この文書集合を、ウェブページの1次検索結果と呼ぶ）を対象として、同様にトピックモデルを適用し、両者のトピック分布を比較する。その結果、図1上半分のように、

ニュース・ブログにおける話題の分布と比較して、検索エンジンAPIを用いることによりニュース・ブログ以外のウェブページから収集された文書集合における話題の分布が相対的に小さいことを示す。

次に、提案手法においては、ニュースもしくはブログのみで観測された話題との関連が強いクエリを用いることにより、ニュース・ブログのみにおいて観測された話題との関連が強いウェブページを収集する（ここで収集された文書集合を、ウェブページの2次検索結果と呼ぶ）。そして、2次検索により収集したウェブページ群に対してトピックモデルを適用することによって推定した話題の分布を、1次検索結果における話題の分布と比較して、提案方式によって収集されたウェブページ群において、より多様性に富んだ話題の分布が観測できることを示す。特に、1次検索結果において、ニュースあるいはブログのみにおいて観測された話題のウェブページだけでなく、ニュースあるいはブログのいずれにおいても観測されなかった多様な話題のウェブページもあわせて収集可能であることを示す。また、[小池13]においては、クエリを手動で選定しウェブページ収集の2次検索を行っていたが、本論文では、クエリの自動選定を行った場合の結果について報告する。結果としては、手動で選定したクエリのうちの約30%のみが共通に選定されたが、収集されたウェブページにおける話題の分布としては、手動選定クエリの場合の約3分の2の話題についてウェブページが収集されるという結果が達成できた。

## 2. 文書収集の手順

本論文で分析対象としたウェブページ、ニュース記事、ブログ記事の収集手順を以下に述べる。本論文で分析対象とした文書数の一覧を表1に示す。

連絡先: 小池 大地, 筑波大学大学院システム情報工学研究科,  
〒305-8573 茨城県つくば市天王台1-1-1, 029-853-5427

表 1: 収集したウェブページ, ニュース, ブログの文書数

ウェブページ	ニュース	ブログ
1 次検索: 593	26,228	20,716
2 次検索 手動: 20,374 自動: 17,247	(朝日: 7,541, 読売: 6,568, 日経: 12,119)	

## 2.1 ウェブページの収集

ウェブページの収集においては, Yahoo! Search BOSS API<sup>\*1</sup> を用い, 検索エンジン API に対してクエリを指定することにより, 日本語のサイトを対象として収集を行った. まず, 初期クエリ  $t_0$  を「東日本大震災」とし, 一度に最大で 1,000 件のウェブページを取得した. その結果, 720 件のウェブページが検索結果として得られた. その後, 主要なニュース・ブログサイトのウェブページを除いたものを分析対象とした. その結果, 1 次検索での分析対象のウェブページは, 593 件となった.

次に, 2 次検索として, 4.2 節の手法により, 手動もしくは自動で選定された語を 2 次クエリ  $t_1$  とし, 各 2 次クエリ, 初期クエリ  $t_0$  との AND 検索によりウェブページの収集を行った. その後, 主要なニュース・ブログサイトのウェブページを除いたものを 2 次検索での分析対象のウェブページとした. 手動で 2 次クエリを選定した場合には, 68,880 件が検索結果として得られ, ニュース・ブログサイトのウェブページを除去した結果の分析対象ウェブページ数は 20,374 件となった. 同様に, 手動で 2 次クエリを選定した場合には, 分析対象ウェブページ数は 17,247 件となった.

## 2.2 ニュース記事の収集

ニュース記事としては, 2011 年 3 月 11 日から 2012 年 7 月 10 日までの日付の記事を, 日経新聞<sup>\*2</sup>, 朝日新聞<sup>\*3</sup>, 読売新聞<sup>\*4</sup> の各新聞社のサイトから収集した 92,772 記事, 63,906 記事, および, 68,239 記事の合計 224,867 記事を用いた. その後, 「東日本大震災」の 1 語がニュース記事中出现するものだけを分析対象とした. その結果, 各新聞社の記事数は, 日経新聞が 12,119 記事, 朝日新聞が 7,541 記事, 読売新聞が 6,568 記事, 合計 26,228 記事となった.

## 2.3 ブログ記事の収集

東日本大震災に関連するブログ記事の収集においては, 東日本大震災との関連性が高い語として, 人手で 26 個の語を選定し, その一つ一つを初期クエリ  $t_0$  とし, ブログ記事を収集した結果を用いた. 初期クエリ  $t_0$  を含む日本語ブログ記事の収集においては, Yahoo! Search BOSS API を利用し, 日本語ブログホスト大手 6 社<sup>\*5</sup> のドメインを対象として, 2012 年 8 月下旬から 9 月上旬に, 2011 年 3 月 11 日以降の日付の記事を対象として, ブログ記事の収集を行った. 検索の際には, 複数のドメインを一度に指定して検索し, 1,000 件の記事を取得する. 次に, ブログ記事検索後, 検索結果の URL をブログサイト単位にまとめる. その結果, 一つの検索クエリあたり約 200 前後のブログサイトが取得される. 次に, 各ブログサイトをドメイン指定し, 初期クエリ  $t_0$  を検索クエリとすることにより, 各ブログサイト中において初期クエリ  $t_0$  を含むブログ

表 2: トピックモデル推定時に指定したトピック数

ウェブページ	ニュース	ブログ
1 次検索: 15	70	60
2 次検索 (手動・自動とも): 90		

記事を収集し, ブログ記事集合を作成する. その後, 「東日本大震災」の 1 語がブログ記事中出现するものだけを分析対象とした. その結果, 分析対象のブログ記事数は, 20,716 記事となった.

## 3. トピックモデルを用いた話題分布の推定

本研究では, トピックモデルとして潜在的ディリクレ配分法 (LDA; Latent Dirichlet Allocation) [Blei03b] を用いる. LDA を用いたトピックモデルの推定においては, 語  $w$  の列によって表現された文書の集合と, トピック数  $K$  を入力として, 各トピック  $z_n$  ( $n = 1, \dots, K$ ) における語  $w$  の確率分布  $P(w|z_n)$  ( $w \in V$ ), 及び, 各文書  $b$  におけるトピック  $z_n$  の確率分布  $P(z_n|b)$  ( $n = 1, \dots, K$ ) を推定する. これらを推定するためのツールとしては, GibbsLDA++<sup>\*6</sup> を用いた. LDA のハイパーパラメータである  $\alpha$ ,  $\beta$  には, GibbsLDA++ の基本設定値である  $\alpha = 50/K$ ,  $\beta = 0.1$  を用いた. LDA ではトピック数  $K$  を人手で与える必要があるが, 本論文では, トピック数を 10 から 100 まで変化させてトピック推定を行い, 得られたトピックを人手で見比べ, トピックの推定結果の性能がより高くなったトピック数を採用するという手順を採った. なお, このツールは推定の際に Gibbs サンプリングを用いているが, その反復回数は 2,000 とした.

## 4. 多様なトピックのウェブページの収集

### 4.1 ニュースおよびブログを利用した多様なトピックの生成

まず, 2.1 節の 1 次検索の手順によって収集したウェブページ (ただし, ニュース・ブログを除く), 2.2 節の手順によって収集したニュース記事, および, 2.3 節の手順によって収集したブログ記事に対して, それぞれ独立にトピックモデルの推定を行った. ただし, 予備実験を経たうえで, それぞれ最も性能よくトピックモデルの推定が行えたトピック数として, 表 2 に示すトピック数を用いた. このうち, ウェブページ集合から推定された 15 トピック, ニュース記事集合から推定された 70 トピック, ブログ記事集合から推定された 60 トピックのうち, 東日本大震災に関連し, かつ, トピックに対応する文書集合において意味的まとまりのあるトピックの数は, それぞれ, 8 トピック, 67 トピック, 40 トピックであった. 次に, これらのトピックのうち, 情報源となった文書がウェブページであるか, ニュース記事であるか, ブログ記事であるかの別を問わず, 同一の話題と考えられるトピックの集約を行ったところ, 合計 62 個の話題に集約された. この 62 個の話題のうち, ウェブページ集合において観測された話題は, 図 2 のベン図のうち「ウェブページ 1 次検索」に示す 7 個の話題である. その他の話題は, ニュースのみでの話題 31 個 (=20+4+7), ブログのみでの話題 12 個 (=9+3), ニュース・ブログの両方での話題 12 個 (=7+5) である.

このように, ウェブページ集合において観測された話題は, 全 62 話題のうちの約 10% 程度に過ぎない. この結果から, ニュー

\*1 <http://developer.yahoo.com/search/boss/>

\*2 <http://www.nikkei.com/>

\*3 <http://www.asahi.com/>

\*4 <http://www.yomiuri.co.jp/>

\*5 fc2.com, yahoo.co.jp, ameblo.jp, goo.ne.jp, livedoor.jp, hatena.ne.jp

\*6 <http://gibbslda.sourceforge.net/>

表 3: ニュース・ブログのみで観測された話題を利用したクエリ拡張におけるクエリ

(a) クエリ拡張: 手動 (全 86 クエリのうち, 太字・下線 26 クエリ は「自動」と共通)

ニュースのみから選定したクエリ (計 49 個)	ブログのみから選定したクエリ (計 22 個)	ニュース・ブログ両方から生成したクエリ (計 15 個)
デザイン, 観光, 公演, 作品, 試合, 避難所, 陛下, 報道, ビール, アンケート, エネルギー, マンション, みずほ, ローン, 営業利益, 沿岸, 学生, 給与, 漁業, 経営, 経団連, 原子炉, 交渉, 交付金, 公的資金, 行方不明者, 高速道路, 国会, 祭り, 市場, 指紋, 事件, 自動車, 需要, 就職, 新幹線, 申請, 水準, 政治, 生徒, 台風, 大阪, 地震保険, 通信, 統一地方選, 百貨店, 貿易収支, 防災, 路線	韓国, 雇用, 子ども, 神社, 地震, 調査, 東証, 東電, 放射線, イベント, フィギュアスケート, メディア, メーカー, 医療, 営業, 小沢一郎, 食品, 脱原発, 中国, 天皇, 日本人, 被災者	義援金, 経済, 原発, 自衛隊, 選手, 増税, 津波, 電力, 福島県, がれき, ボランティア, 家族, 住宅, 世界, 政府

(b) クエリ拡張: 自動 (全 82 クエリのうち, 太字・下線 26 クエリ は「手動」と共通)

ニュースのみから選定したクエリ (計 49 個)	ブログのみから選定したクエリ (計 20 個)	ニュース・ブログ両方から生成したクエリ (計 13 個)
デザイン, 観光, 公演, 作品, 試合, 避難所, 陛下, 経済, 調査, カ月, システム, ドル, ポイント, 以上, 価格, 過去, 会議, 開発, 学校, 企業, 記者会見, 宮城, 宮城県, 建設, 国際, 今年, 作業, 資金, 受け, 情報, 人口, 生産, 対策, 対象, 大学, 男性, 知事, 地域, 東京, 東日本, 年度, 派遣, 売上高, 発表, 販売, 被害, 被災, 払い, 問題	韓国, 子ども, 東証, 東電, 放射線, 報道, 神社, 義援金, 自衛隊, 増税, 津波, 活動, 検査, 工場, 国民, 紹介, 処理, 対応, 病院, ライブ	原発, 選手, 電力, 福島県, 雇用, 地震, 午後, 自分, 写真, 首相, 日本, 被災地, 復興

ス・ブログ以外のウェブページを対象として, 検索エンジン API の上位 1,000 位以内の範囲において, ニュース・ブログにおいて観測されるような多様な話題の文書を収集することが容易でないことが判明した.

#### 4.2 各トピックに対応するウェブページの収集

次に, ニュース記事のみにおいて観測された 31 話題, ブログ記事のみにおいて観測された 12 話題, ニュース記事とブログ記事の両方で観測され, ニュース・ブログ以外のウェブページでは観測されなかった 12 話題に対応する各トピック  $z_n$  において, 確率値  $P(w|z_n)$  が上位 20 位以内となる語  $w$  のうち, 初期クエリ  $t_0$  (=「東日本大震災」) との AND 検索を行う 2 次クエリ  $t_1$  を手動, および, 自動で選定し, 検索エンジン API を用いたウェブページ検索により, 各 2 次クエリごとに最大で 1,000 ページを取得した. 2 次クエリ  $t_1$  を手動で選定する際には, 各トピック  $z_n$  における話題を適切に反映する語を 1 語選定した\*7. 一方, 2 次クエリ  $t_1$  を自動で選定する際には, 確率値  $P(w|z_n)$  が最大となる語  $w$  を選定した. 手動および自動で選定されたクエリを表 3 に示す. 両者の間で共通となったクエリは, 表中で 太字・下線 で示す 26 個 (約 30%) であった. ただし, 自動で選定されたクエリのうち, 手動で選定されたクエリと一致しないものの中にも, 各トピックの表す話題との関連性が高いものは多数存在している.

#### 4.3 分析

次に, 2 次検索によって収集されたウェブページを対象として, トピックモデルの推定を行った. ただし, トピック数をと

しては, 予備実験を経たうえで, 最も性能よくトピックモデルの推定が行えたトピック数として, 表 2 に示すトピック数 90 を用いた. このうち, 東日本大震災に関連し, かつ, トピックに対応する文書集合において意味的まとまりのあるトピックの数は, クエリを手動で選定した場合は, 50 トピック, クエリを自動で選定した場合は, 44 トピックであった. であった. さらに, 前節において人手で集約した 62 個の話題に対して, 同様に, 50 トピック, および, 44 トピックとの対応付けを行った. その結果, 図 2 のベン図に示すように, 50 トピックに対しては新たに 9 個の話題が追加され, ウェブページ集合に対する話題の数は 36 個となった. また, 44 トピックに対しては新たに 4 個の話題が追加され, ウェブページ集合に対する話題の数は 24 個となった. この結果から分かるように, 自動選定されたクエリを用いることにより, 手動で選定されたクエリにより収集したウェブページ集合における 36 個の話題のうち約 3 分の 2 に当たる 22 個の話題を復元することができ, さらに, 新規の話題 2 個に相当するウェブページを収集することができた.

この結果のうち, 特に, 自動選定されたクエリを用いた 2 次検索によって収集されたウェブページにおける話題の分布に注目すると, 1 次検索結果において観測された 6 個の話題を除く 18 個の話題に対応するウェブページを新規に収集できたことが分かる. 以上の結果から, 検索エンジン API を用いた 1 次検索においては, ニュース・ブログ以外のウェブページからは収集困難であった多様な話題が, 提案手法を用いることにより収集可能となることが示された.

## 5. 関連研究

本論文に関連して, Web ページの検索結果を分類し, 各分類に対して適切な要約文を付与するという手法 [原島 10], および, 検索された個々の Web ページに対してラベルの付与を行い, 付与されたラベルに基づいて分類を行う手法 [戸田 05, de Winter 07, 馬場 09], 階層的なトピックの体系を推定する手

\*7 ここで, 同一の話題に相当するトピックが複数存在する可能性があるが, 本論文において, 各トピックから 2 次クエリを手動で選定する際には, 他のトピックにおけるクエリ候補を参照せず, 独立に 2 次クエリを選定した. このように, 本論文における評価結果は, 2 次クエリを手動で選定することによって, 話題の広がりか最も大きくなる上限値を示したものと位置付けることができる. ただし, 手動クエリ選定, 自動クエリ選定, いずれの場合も, 各トピックからは, 重複する 2 次クエリを選定しないという制約を課すことにより, 話題の広がりをより大きくできると考えられる.

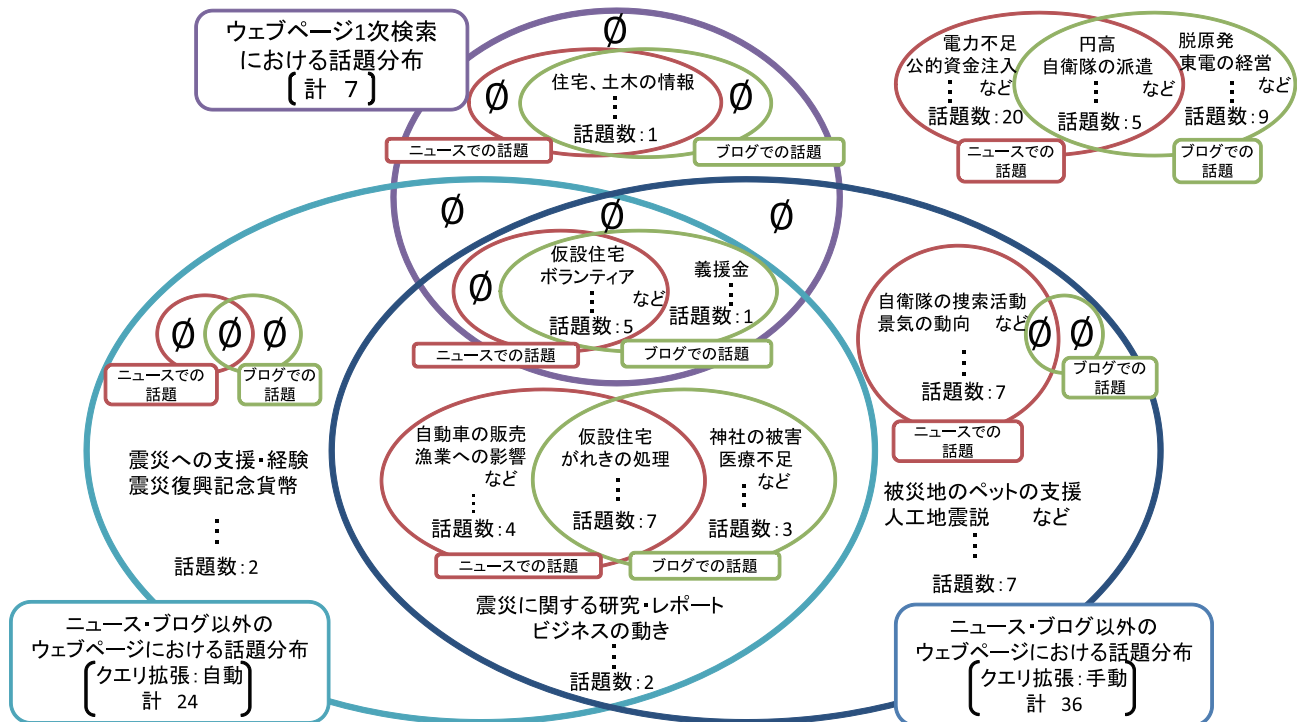


図 2: ウェブページ・ニュース・ブログにおける話題の分布

法 [Blei03a] 等が提案されている。これらの手法においては、いずれも、閲覧対象の文書集合のみを用いて、ファセット体系およびファセットラベルに相当する情報を抽出している。また、メタ検索エンジンにおいてウェブページ検索結果の上位 200 記事程度を対象にして、検索結果のクラスタリングおよびラベル付けをした結果を提示するサービスとして、Yippy\*<sup>8</sup> が知られている。一方、本論文では、「検索エンジン API に対してクエリを与えることにより、ウェブページを収集する」というタスクを設定し、クエリを変更して複数回 API アクセスを行うことを許容し、かつ、ニュース・ブログといった外部言語資源も援用するという枠組みのもとで、できるだけ多様な話題のウェブページを収集する方式を提案している。

また、本論文に関連して、TREC の Web Track の diversity タスク [Clarke12] や NTCIR の INTENT タスク [Song11] における文書ランキングタスクにおいては、ウェブ検索結果においてできるだけ多様な話題のウェブページを上位に順位付けすることを要求する仕様のもとで評価型タスクを行っている。同様に、NTCIR の INTENT タスク [Song11] におけるサブトピックマイニングタスクにおいては、クエリについてのサブトピックを列挙する課題を設定し、評価型タスクを行っている。

## 6. おわりに

本論文では、検索エンジン API を用いたウェブページ収集タスクにおいて、ニュースもしくはブログにおいてのみ観測され、ニュースあるいはブログ以外のウェブページからは、検索エンジン API によって収集されなかった話題の文書を選択的に収集することにより、検索エンジン API を用いてできる限り多様な話題のウェブページを収集する方式を提案した。特に、本論文では、ウェブページの 2 次検索の際のクエリを自動選定した場合の結果を報告したが、今後、ウェブページの 2 次検索結果においてより多様な話題のウェブページを収集するため

に、クエリ自動選定方式を改善することが重要な課題である。

また、現在、各話題について、2 次検索によって収集されたウェブページ中の記載内容と、ニュース・ブログから収集された記事における記載内容の比較対照分析作業を行っている。この比較対照分析作業によって、ニュース・ブログでは取り上げられておらず、ウェブページ固有の情報とみなせる記載内容が、2 次検索によって収集されたウェブページ中にどの程度含まれるのかの検証を行う予定である。この検証を通して、提案手法にしたがって検索エンジン API を用いることにより、ニュース・ブログ以外のウェブページから、従来アクセスが困難であった情報が収集できるか否かが明らかになると考えている。

## 参考文献

- [馬場 09] 馬場康夫, 黒橋禎夫: キーワード蒸留型クラスタリングによる大規模ウェブ情報の俯瞰, 情報処理学会論文誌, Vol. 50, No. 4, pp. 1399–1409 (2009).
- [Blei03a] Blei, D. M., Griffiths, T. L., Jordan, M. I. and Tenenbaum, J. B.: Hierarchical Topic Models and the Nested Chinese Restaurant Process, *NIPS'03* (2003).
- [Blei03b] Blei, D. M., Ng, A. Y. and Jordan, M. I.: Latent Dirichlet Allocation, *Journal of Machine Learning Research*, Vol. 3, pp. 993–1022 (2003).
- [Clarke12] Clarke, C. L. A., et al.: Overview of the TREC-2012 Web Track, *Proc. TREC-2012* (2012).
- [de Winter07] de Winter, W. and de Rijke, M.: Identifying Facets in Query-Biased Sets of Blog Posts, *Proc. ICWSM*, pp. 251–254 (2007).
- [原島 10] 原島純, 黒橋禎夫: PLSI を用いたウェブ検索結果の要約, 言語処理学会第 16 回年次大会論文集, pp. 118–121 (2010).
- [小池 13] 小池大地, 牧田健作, 宇津呂武仁, 河田容英, 吉岡真治, 神門典子: 検索エンジン API を用いたウェブページ収集におけるトピックの多様性, 第 5 回 DEIM フォーラム論文集 (2013).
- [Song11] Song, R., et al.: Overview of the NTCIR-9 INTENT Task, *Proc. 9th NTCIR Workshop Meeting*, pp. 82–105 (2011).
- [戸田 05] 戸田浩之, 中渡瀬秀一, 片岡良治: 特徴的な固有表現を用いたラベル指向ナビゲーション手法の提案, 情報処理学会論文誌: データベース, Vol. 46, No. SIG 13(TOD 27), pp. 40–52 (2005).

\*8 <http://yippy.com/>