



HAL
open science

Evaluating search engines and defining a consensus implementation

Ahmed Kamoun, Patrick Maillé, Bruno Tuffin

► **To cite this version:**

Ahmed Kamoun, Patrick Maillé, Bruno Tuffin. Evaluating search engines and defining a consensus implementation. VALUETOOLS 2019 - 12th EAI International Conference on Performance Evaluation Methodologies and Tools, Mar 2019, Palma de Majorque, Spain. pp.1-10. hal-01852650

HAL Id: hal-01852650

<https://inria.hal.science/hal-01852650>

Submitted on 2 Aug 2018

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Evaluating search engines and defining a consensus implementation

Ahmed Kamoun
IMT Atlantique
Rennes, France
ahmed.kamoun@imt-atlantique.net

Patrick Maillé
IMT Atlantique
Rennes, France
patrick.maille@imt.fr

Bruno Tuffin
Inria Rennes Bretagne Atlantique
Campus Universitaire de Beaulieu, 35042 Rennes Cedex, France
bruno.tuffin@inria.fr

August 2, 2018

Abstract

Different search engines provide different outputs for the same keyword. This may be due to different definitions of relevance, and/or to different knowledge/anticipation of users' preferences, but rankings are also suspected to be biased towards own content, which may be prejudicial to other content providers. In this paper, we make some initial steps toward a rigorous comparison and analysis of search engines, by proposing a definition for a consensual relevance of a page with respect to a keyword, from a set of search engines. More specifically, we look at the results of several search engines for a sample of keywords, and define for each keyword the visibility of a page based on its ranking over all search engines. This allows to define a score of the search engine for a keyword, and then its average score over all keywords. Based on the pages visibility, we can also define the consensus search engine as the one showing the most visible results for each keyword. We have implemented this model and present an analysis of the results.

1 Introduction

Search Engines (SEs) play a crucial role in the current Internet world. If you wish to reach some content, except if you have a specific target in mind, you dial keywords on an SE through a web browser to discover the (expected) most relevant content. The number of searches worldwide per year is not precisely known, but just talking about Google, it is thought that they handle at least two trillions of requests per year, and that it can even be much more than that¹. As a consequence, if you are a small content provider or a new comer, your visibility and business success will highly depend on your ranking on SEs.

SEs are regularly accused of biasing their rankings² by not only trying to provide as an output an ordered list of links based on *relevance*, but to also include criteria based on revenues it could drive. The problem was brought in 2009 by Adam Raff, co-founder of the price-comparison company Foundem, saying that Google was voluntarily penalizing his company in rankings with respect to Google's own services. Such a behavior would indeed be rational from the SE perspective, as it could yield significant revenue increases; a mathematical model highlighting the optimal

¹<https://searchengineland.com/google-now-handles-2-999-trillion-searches-per-year-250247>

²See for example

<https://www.technologyreview.com/s/610275/meet-the-woman-who-searches-out-search-engines-bias-against-women-and-minorities/>

non-neutral—i.e., not based on relevance only—strategies of SEs is for example described in [3]. The issue led to the use of expression *search neutrality debate*, in relation to the *net neutrality debate* where Internet Service Providers are accused of differentiating service at the packet level to favor some applications, content, or users. Indeed, similarly, new valid content can hardly be reached if not properly considered by SEs. This is now an important debate worldwide [2, 5, 8]. But while defining a neutral behavior of ISPs at the network level is quite easy, a neutral behavior for SEs involves having a clear definition of relevance. Up to now this relevance is defined by SE-specific algorithms such as PageRank [7], that can additionally be (and are claimed to be) refined by taking into account location, cookies, etc. The exact used algorithms and their relation to relevance are sometimes hard to know without requiring a total transparency of SEs and free access to their algorithms, which they are reluctant to disclose.

Because of the different used algorithms, it is interesting to compare search engines, for example by giving them a grade (or score). It is often said that if someone is not happy, she can just switch, she is just one click away from another SE. But while it could be true in a fully competitive market, it is not so easy in practice with SEs since most people just know one or two SEs and do not have a sufficient expertise to evaluate them and switch. As of May 2018, Statcounter Global Stats³ gives worldwide a market share of 90.14% to Google, 3.24% to Bing, 2.18% to Baidu, 2.08% to Yahoo!...

Our paper has several goals:

- First, to propose a so-called *consensus SE*, defined such as some “average” behavior of SEs, based on the idea that this average SE should be closer to one truly based on relevance. Considering a list of several SEs, we give a score to all the provided links, by weighing them by the position click-through-rate on each SE, estimating the position-dependent probability to be clicked. The consensus SE then ranks links according to their score.
- To give a score to SEs, comparing their results to the consensus SE. From the score of links, we can give a score to SEs by summing the scores of the presented lists weighted by their positions. It then allows us to rank the SEs themselves and show which one seems the closest to the “optimal” consensus SE, for a single keyword and for a list of keywords.
- To discuss and compare the behavior of SEs with respect to requests in practice. We have implemented and tested our model, computing grades for SEs and distributions of scores in terms of requests. From the rankings of SEs for any keyword, we can also investigate if there is a suspect deviation of some SEs toward their own content with respect to competitors. This would help to detect violations to a (potential) search neutrality principle.

Note that algorithms comparing rankings exist in the literature, see [6] and references therein, based on differences between vectors, but there is to our knowledge no algorithm like ours taking into account the click-through-rates (weights) associated to positions.

The rest of the paper is organized as follows. Section 2 presents the model: the definition of link scores for any keyword, the corresponding SE score as well as the consensus SE. Section 3 presents an implementation of this model in practice and compares the most notable SEs. Finally, Section 4 concludes this preliminary work and introduces the perspectives of extension.

2 Scoring model and consensus search engine

We consider n SEs, m keywords representative of real searches, and a finite set of ℓ pages/links corresponding to all the results displayed for the whole set of searches. When dialing a keyword, SEs rank links. We will limit ourselves to the first displayed page of each SE, considered here for simplicity the same number a for all SEs, but we could consider different values for each SE, and even $a = \ell$.

³<http://gs.statcounter.com/search-engine-market-share>

The ranking is made according to a score assigned to each page for the considered keyword. This score is supposed to correspond the relevance of the page. According to their rank, pages are more or less likely to be seen and clicked. The probability to be clicked is called the click-through-rate (CTR) [4]; it is in general SE-, position- and link- dependent, but we assume here for convenience, and as commonly adopted in the literature, a separability property: the CTR of link i at position l is the product $q'_i q_l$ of two factors, q'_i depending on the link i only, and q_l depending on the position l only. We typically have $q_1 \geq q_2 \geq \dots \geq q_a$. The difficulty is that the link-relative term q'_i , upon which a “neutral” ranking would be based, is unknown. But the position-relative terms q_l can be estimated, and we assume in this paper that they are known and the same on all SEs, i.e., that SEs’ presentation does not influence the CTR. We then make the following reasoning:

- SEs favor the most relevant (according to them) links by ranking them high, that is, providing them with a good position and hence a high *visibility*, which can be quantified by the position-relative term q_l ;
- for a given keyword, a link that is given a high visibility by all SEs is likely to be “objectively relevant”. Hence we use the average visibility offered by the set of SEs to a link, as an indication of the link relevance for that keyword. We will call that value the *score* of the link for the keyword, which includes the link-relative term q'_i .
- We expect that considering several SEs will average out the possible biases introduced by individual SEs, when estimating relevance; also, the analysis may highlight some SEs that significantly differ from the consensus for some sensitive keywords, which would help us detect non-neutral behaviors.

2.1 Page score

The notion of score of the page as defined by SEs is (or should) be related to the notion of relevance for any keyword. As briefly explained in the introduction, the idea of relevance is subjective and depends on so many possible parameters that it can hardly be argued that SEs do not consider a valid definition without knowing the algorithm they use. But transparency is very unlikely because the algorithm is the key element of their business.

In this paper, as explained above we use a different and original option for defining the score, as the exposition (or visibility) provided by all SEs, which can be easily computed.

Formally, for page i and keyword k , the score is the average visibility over all considered SEs:

$$R_{i,k} := \frac{1}{n} \sum_{j=1}^n q_{\pi_j(i,k)} \quad (1)$$

where $\pi_j(i, k)$ denotes the position of page i on SE j for keyword k . In this definition, if a page is not displayed by an SE, the CTR is taken as 0. Another way to say it is to define a position $a + 1$ for non displayed pages, with $q_{a+1} = 0$.

2.2 Search engine score

Using the score of pages (corresponding to their importance), we can define the score of an SE j for a given keyword k as the total “page score visibility” of its results for that keyword. Mathematically, that SE score $S_{j,k}$ can be expressed as

$$S_{j,k} := \sum_{\text{pages } i} q_{\pi_j(i,k)} R_{i,k},$$

where again $q_p = 0$ if a page is ranked at position $p \geq a + 1$ (i.e., not shown), and for each page i , $R_{i,k}$ is computed as in Eq. (1).

The SE score can also be computed more simply, by just summing on the displayed pages:

$$S_{j,k} = \sum_{p=1}^a q_p R_{\tilde{\pi}_j(p,k),k} \quad (2)$$

where $\tilde{\pi}_j(p,k)$ is the page ranked at the p^{th} position by SE j for keyword k , i.e., $\tilde{\pi}_j(\cdot,k)$ is the inverse permutation of $\pi_j(\cdot,k)$.

The higher an SE ranks highly exposed pages, the higher its score. The score therefore corresponds to the exposition of pages that are well-exposed on average by SEs.

To define the score of SE j , for the whole set of keywords, we average over all keywords:

$$S_j := \frac{1}{m} \sum_{k=1}^m S_{j,k}. \quad (3)$$

2.3 Consensus search engine

From our definitions of scores in the previous subsection, we can define the *consensus SE* as the one maximizing the SE score for each keyword. Formally, for a given keyword k , the goal of the consensus SE is to find an ordered list of the ℓ pages (actually, getting the first a is sufficient), where $\pi^{(k)}(p)$ is for the page at position p , such that

$$\pi^{(k)}(\cdot) = \operatorname{argmax}_{\pi(\cdot)} \sum_{p=1}^a q_p R_{\pi(p),k}.$$

Note that this maximization is easy to solve: it suffices to order the pages such that $R_{\pi^{(k)}(1),k} \geq R_{\pi^{(k)}(2),k} \geq \dots$, i.e., to display pages in the decreasing order of their score (visibility).

The total score of the consensus SE can then also be computed, and is straightforwardly maximal.

3 Analysis in practice

We have implemented in Python a web crawler that looks, for a set of keywords, the results provided by nine different search engines. From those results, the scores can be computed as described in the previous section, as well as the results and score of a consensus SE. The brute results can be found at <https://partage.mines-telecom.fr/index.php/s/aG3SYhVYPtRCBKH>. The code to get page URLs is adapted to each SE, because they display the results differently. It also deals with results display that can group pages and subpages (that is, lower level pages on a same web site) that could be treated as different otherwise. Another solved issue is that *a priori* different URLs can lead to the same page. It is for example the case of <http://www.maps.com/FunFacts.aspx>, <http://www.maps.com/funfacts.aspx>, <http://www.maps.com/FunFacts.aspx?nav=FF>, <http://www.maps.com/FunFacts>, etc. It can be checked that they actually lead to the same web page output when accessing the links proposed by the SEs. Note on the other hand that it requires a longer time for our crawler to get to each page and check whether the URL gets modified.

We (arbitrarily) consider the following set of nine SEs among the most popular, in terms of number of requests according to <https://www.alexa.com/siteinfo>:

- Google
- Yahoo!
- Bing
- AOL
- ask.com

- duckduckgo
- Ecosia
- StartPage
- Qwant.

We include SEs such as Qwant or StartPage, which are said to respect privacy and neutrality. We also clear the cookies to prevent them from affecting the results (most SEs use cookies to learn our preferences).

We consider 216 different queries included in February 2018 common searches. The choice is based on the so-called trending searches in various domains according to <https://trends.google.fr/trends/topcharts>. We arbitrarily chose keywords in different categories to cover a large range of possibilities.

We limit ourselves to the first page of search engines results, usually made of the first 10 links. That is, we let $a = 10$. The click-through rates q_p are set as measured in [1] and displayed in Table 1.

q_1	q_2	q_3	q_4	q_5	q_6	q_7	q_8	q_9	q_{10}
0.364	0.125	0.095	0.079	0.061	0.041	0.038	0.035	0.03	0.022

Table 1: CTR values used in the simulations, taken from [1]

3.1 Search engines scores

Table 2 provides the average scores of the nine considered search engines, as well as that of the consensus SE, according to Eq. (3). We also include the 95% confidence intervals that would be obtained (abusively) assuming requests are independently drawn from a distribution on all possible requests. Under the same assumption, we can also implement statistical tests to determine

SE	Score
Google	0.0832 ± 0.0045
Yahoo	0.1103 ± 0.0030
Bing	0.0933 ± 0.0045
AOL	0.1055 ± 0.0036
Ask	0.0211 ± 0.0006
DuckDuckGo	0.1106 ± 0.0029
Ecosia	0.1071 ± 0.0033
StartPage	0.0816 ± 0.0046
Qwant	0.0906 ± 0.0048
Consensus	0.1332 ± 0.0026

Table 2: SE scores, and 95% confidence intervals half-widths.

whether the scores of search engines are significantly different. The corresponding p -values are given in Table 3. For two search engines, the p -value is the probability of error when rejecting the hypothesis they have similar scores. A small value indicates a statistically significant difference between both search engines (1% means 1% chance of error).

We can remark a group of four SEs with scores above the others: DuckDuckGo, Yahoo!, Ecosia, and AOL, around 0.11. The statistical analysis using the p -value allows to differentiate even more, with DuckDuckGo and Yahoo! as a first group, and Ecosia, and AOL slightly below. Then, Bing and Qwant get scores around 0.09 (and can not be strongly differentiated from the p -value in

	Yahoo	Bing	AOL	Ask	DuckDuckGo	Ecosia	StartPage	Qwant	Consensus
Google	1.3e-38	5.5e-05	9.7e-24	7.7e-70	6.3e-42	5.4e-30	1.3e-01	5.4e-03	6.8e-82
Yahoo		5.0e-23	3.5e-08	1.5e-131	5.4e-01	1.9e-04	1.0e-40	2.4e-25	2.4e-129
Bing			5.8e-11	2.6e-81	3.1e-23	2.8e-15	6.6e-06	6.1e-02	4.0e-70
AOL				6.1e-112	3.9e-07	1.3e-01	2.0e-26	1.9e-13	1.6e-78
Ask					4.5e-135	5.1e-120	4.5e-67	8.3e-75	4.0e-163
DuckDuckGo						4.5e-05	4.5e-42	4.4e-25	1.4e-130
Ecosia							5.6e-32	1.8e-17	2.0e-91
StartPage								9.4e-04	2.3e-82
Qwant									1.6e-67

Table 3: p -values for the tests comparing the average scores of search engines (T-test on related samples of scores).

Table 3)), and Google and StartPage around 0.082 (since StartPage is based on Google, close results were expected). Finally, quite far from the others, Ask.com has a score around 0.02.

The consensus SE has a score of 0.133, significantly above all the SEs as shown in Table 3.

3.2 Analysis

Figure 1 displays by SE the percentage of common results with the consensus SE for each position range. Again for all our searched keywords, we count the proportion of links in the 1st position correspond to the link in 1st position in the consensus SE, then do the same for the first 2 positions, then for the first 3, etc.

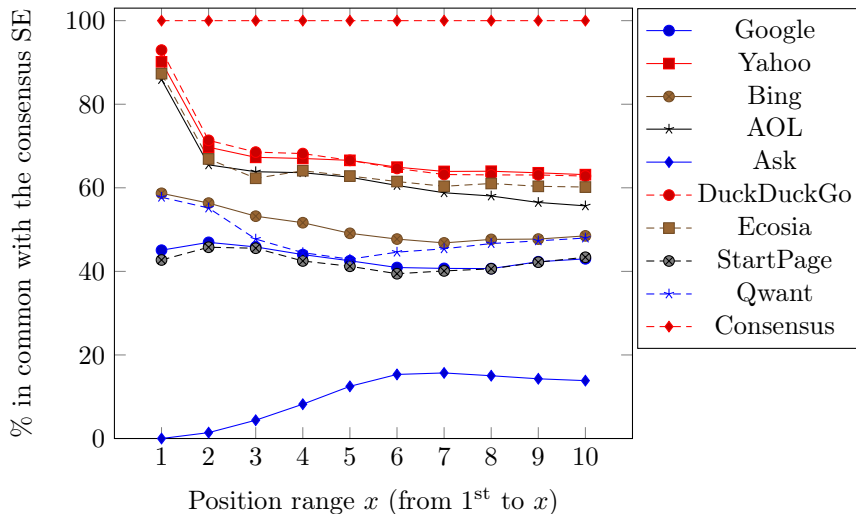


Figure 1: Similarities in position with the consensus

The results are consistent with the previous tables: the figure highlights the same groups of SE, while Ask.com clearly is far from the consensus.

We also draw in Figure 2 the distribution of the score of SEs relatively to the consensus where on the x -axis, we have the pages ordered (for each SE) by the relative score from the largest to the smallest.

It allows to see the variations of score per SE. Again the same SE groups appear, but the information is stronger than just the mean. For the first quarter of requests, scores are close for all

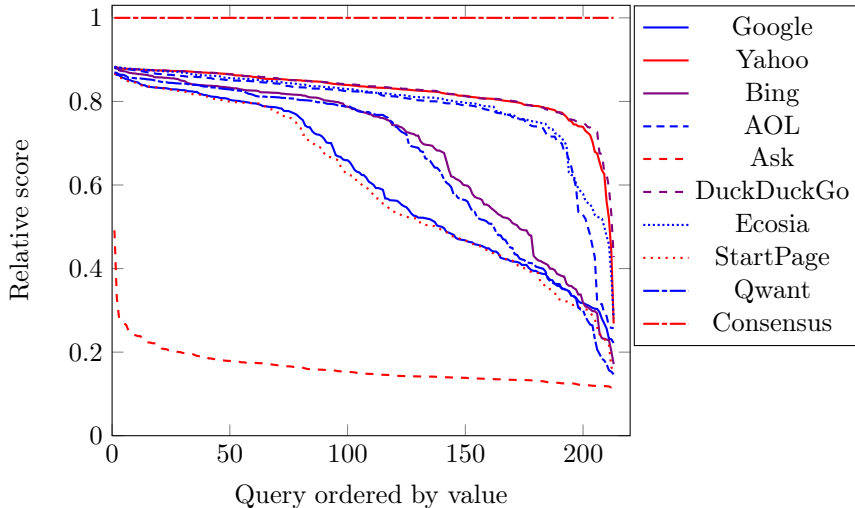


Figure 2: Distribution of scores relative to the consensus, from largest to smallest.

SEs except Ask.com, the difference becomes significant later with some SEs which cannot keep up with the best ones.

To identify deviations from other search engines, we highlight respectively in Tables 4 and 5 for each SE the (ordered) 10 queries with the highest and lowest relative score with respect to the consensus SE. Those queries correspond to the extreme left (for Table 4) and extreme right (for Table 5) of Fig. 2.

Search terms displayed in Table 4 appear to be quite complex-or specific-searches, for which there is not much room for disagreement among SEs.

On the other hand, Table 5 shows, for each SE, the terms for which they most disagree with the consensus, which may help highlight non-neutral behaviors. For example, it is interesting to note that Bing, the Microsoft-operated SE, differs most from the consensus (hence, from the other SEs) on some sensitive searches such as `skype`, `gmail`, `youtube`, and `facebook`. Similarly, AOL strongly disagrees with the consensus for `yahoomail`, `mail`, and `messenger`. While Qwant gets low scores for `google maps`, `msn`, `outlook`, `news`, `google drive`, `skype`, and `gmail`, all involving SE-controlled services. Finally, let us note that we may also have a look at searches like `restaurant` or `cnn`, for which Google is far from the consensus: is that to favor its own news and table-booking services? Our study is too preliminary to draw definite conclusions, but can help raise such questions.

4 Conclusions

In this paper, we have defined a measure of relevance of web pages for given queries based on the visibility/response from a whole set of search engines. This relevance takes into account the position of the web page thanks to a weight corresponding to the click-through-rate of the position. It then allowed to define a score of a search engine for a given query, and the average score for a whole set of queries.

We designed a tool in Python allowing to study the scores of nine known search engines and to build the consensus ranking maximizing the SE score, for a set of more than two hundred queries. A first analysis suggests that there are significant differences among search engines, that may help identify some sensitive terms subject to biases in the rankings.

We finally note that our method does not provide an absolute-value score for each SE allowing to decide which is the best one, but rather indicates whether an SE agrees with the others. The user may very well prefer an SE that is far from our consensus ranking, especially if that SE better

takes her preferences into account when performing the ranking.

In a follow up of this preliminary work, we plan to design statistical tests of potentially intentional deviations by search engines from their regular behavior, to highlight if non-neutrality by search engines can be detected and harm some content providers. This would be a useful tool within the search neutrality debate.

References

- [1] R. Dejarnette. Click-through rate of top 10 search results in Google, 2012. <http://www.internetmarketingninjas.com/blog/search-engine-optimization/click-through-rate>, last accessed June 68, 2017.
- [2] Inria. Inria’s response to ARCEP consultation about network neutrality, 2012.
- [3] P. L’Ecuyer, P. Maillé, N. Stier-Moses, and B. Tuffin. Revenue-maximizing rankings for online platforms with quality-sensitive consumers. *Operations Research*, 65(2):408–423, 2017.
- [4] P. Maillé, E. Markakis, M. Naldi, G. Stamoulis, and B. Tuffin. An overview of research on sponsored search auctions. *Electronic Commerce Research Journal*, 12(3):265–300, 2012.
- [5] P. Maillé and B. Tuffin. *Telecommunication Network Economics: From Theory to Applications*. Cambridge University Press, 2014.
- [6] A. Mowshowitz and A. Kawaguchi. Measuring search engine bias. *Information Processing & Management*, 41(5):1193 – 1205, 2005.
- [7] M. Page, S. Brin, R. Motwani, and T. Winograd. The pagerank citation ranking: Bringing order to the web. Technical Report 1999-66, Stanford InfoLab, November 1999. Previous number = SIDL-WP-1999-0120.
- [8] J. D. Wright. Defining and measuring search bias: Some preliminary evidence. George Mason Law & Economics Research Paper 12-14, George Mason University School of Law, 2012.

Google	Yahoo	Bing	AOL	Ask	DuckDuckGo	Ecosia	StartPage	Qwant
(0.8667) how many days until christmas	(0.8799) how to cook quinoa	(0.8835) how to cook quinoa	(0.8799) how to cook quinoa	(0.4915) what does hmu mean	(0.8821) how to take a screenshot on a mac	(0.8821) how to take a screenshot on a mac	(0.8667) how many days until christmas	(0.8696) cricbuzz
(0.8635) how much house can i afford	(0.8793) how to take a screenshot on a mac	(0.8815) how much house can i afford	(0.8793) how to take a screenshot on a mac	(0.3301) how to draw a doghow to get rid of blackheads	(0.8813) how many days until christmas	(0.8818) MercadoLibre	(0.8635) how much house can i afford	(0.8681) how much house can i afford
(0.8626) how many days till christmas	(0.8792) what time is sunset	(0.8742) how to take a screenshot on a mac	(0.8792) what time is sunset	(0.2802) craigslist	(0.8799) how to cook quinoa	(0.8813) how many days until christmas	(0.8568) how many days till christmas	(0.8646) ebay kleinanzeigen
(0.8581) omegle	(0.8782) speedometer test	(0.8728) how to take a screenshot	(0.879) MercadoLibre	(0.2756) who sings this song	(0.8796) crikbuzz	(0.8806) how to screenshot on mac	(0.8557) what is my ip address	(0.8638) how to write a cover letter
(0.8559) what is my ip address	(0.8782) what is my ip	(0.8724) how to download videos from youtube	(0.8788) crikbuzz	(0.2694) how to make french toast	(0.8792) what time is sunset	(0.8799) flipkart	(0.8513) home- depot	(0.8633) crikbuzz
(0.8531) national basketball association	(0.878) cricbuzz	(0.8713) how many centimeters in an inch	(0.8783) what is my ip	(0.2649) restaurant	(0.8782) flipkart	(0.879) what time is sunset	(0.85) what is my ip	(0.8605) how many ounces in a liter
(0.852) when we were young	(0.8779) national basketball association	(0.8688) what time is it in cal- ifornia	(0.8766) how many mb in a gb	(0.2509) tiempos	(0.8782) irctc	(0.8788) weather	(0.8479) national basketball association	(0.8564) how to take a screenshot on a mac
(0.85) what is my ip	(0.8773) weather	(0.8681) what time is it in lon- don	(0.8762) how to write a check	(0.2442) amazon	(0.8782) what is my ip	(0.8774) how many days till christmas	(0.8466) how old is justin bieber	(0.8564) juegos
(0.8479) euro 2016	(0.8771) what time is it in lon- don	(0.867) why is the sky blue	(0.8751) what is my ip address	(0.2421) mailen	(0.8774) how much house can i afford	(0.8774) how to cook quinoa	(0.8461) how to take a screenshot on a mac	(0.8561) what time is it in lon- don
(0.8466) how many people are in the world	(0.8768) why is the sky blue	(0.8668) crikbuzz	(0.8745) weather	(0.2402) national basketball association	(0.8772) speedometer test	(0.877) tubemate	(0.8449) how many people are in the world	(0.8553) bed 365

Table 4: Per SE, ordered list of 10 queries with the largest relative score with respect to the consensus (and their relative scores)

Google	Yahoo	Bing	AOL	Ask	DuckDuckGo	Ecosia	StartPage	Qwant
(0.2219) convertidos	(0.2686) yahoomail	(0.1713) skype	(0.2576) when is fathers day	(0.115) how to take a screenshot on a mac	(0.4272) traduttur	(0.2803) bbc news	(0.1533) beeg	(0.1469) google maps
(0.2309) how to draw a doghow to get rid of blackheads	(0.44) mail	(0.2025) ikea	(0.2581) yahoomail	(0.1153) what is my ip address	(0.5306) how many ounces in a quart	(0.3822) where are you now	(0.1647) convertidos	(0.1533) minecraft
(0.2322) restaurant	(0.5195) how to make love	(0.2287) gmail	(0.2609) traductor google	(0.1179) how to screenshot on mac	(0.5929) how to start a business	(0.4572) how to make money	(0.2158) daily mail	(0.1564) msn
(0.2564) how many ounces in a quart	(0.5607) who sings this song	(0.2301) youtube	(0.2862) when is mothers day	(0.1183) football association	(0.6261) messenger	(0.4779) games	(0.2213) omegle	(0.1712) news
(0.2686) cnn	(0.6264) messenger	(0.2318) youtube mp3	(0.2936) where are you now	(0.1188) juegos	(0.6319) hotmail	(0.5005) how tall is kevin hart	(0.2334) restaurant	(0.1801) google drive
(0.28) ryanair	(0.6351) oranges	(0.2332) pokemon go	(0.3147) mail	(0.119) how many centimeters in an inch	(0.6739) who sings this song	(0.5194) tiempos	(0.2416) how to make love	(0.1844) outlook
(0.2957) putlocker	(0.6552) how do you spell	(0.2452) ryanair	(0.3156) messenger	(0.1192) irctc	(0.6757) euro 2016	(0.5219) myn	(0.2424) ryanair	(0.1957) skype
(0.3033) what time is it in australia	(0.6756) euro 2016	(0.2928) aleg	(0.3196) what is your name	(0.1193) how much house can i afford	(0.7268) zalando	(0.5238) how to make money fast	(0.2591) mincraft	(0.2063) zara
(0.3038) instagram	(0.6794) what is the temperature	(0.2981) facebook	(0.427) euro 2016	(0.1193) weather	(0.7293) what is the temperature	(0.528) how old is hillary clinton	(0.2791) cnn	(0.2318) youtube mp3
(0.3059) traduttore	(0.7037) how many weeks in a year	(0.2981) pandora	(0.4532) who sings this song	(0.1196) how many people are in the world	(0.7296) how to make pancakes	(0.5438) when we were young	(0.2881) how to draw a doghow to get rid of blackheads	(0.2641) gmail

Table 5: Per SE, ordered list of 10 queries with the smallest relative score with respect to the consensus (and their relative score).