

TREC 2015 Total Recall Track Overview

Adam Roegiest, *University of Waterloo*
Gordon V. Cormack, *University of Waterloo*
Maura R. Grossman, *Wachtell, Lipton, Rosen & Katz**
Charles L.A. Clarke, *University of Waterloo*

1 Summary

The primary purpose of the Total Recall Track is to evaluate, through controlled simulation, methods designed to achieve very high recall – as close as practicable to 100% – with a human assessor in the loop. Motivating applications include, among others, electronic discovery in legal proceedings [2], systematic review in evidence-based medicine [11], and the creation of fully labeled test collections for information retrieval (“IR”) evaluation [8]. A secondary, but no less important, purpose is to develop a sandboxed virtual test environment within which IR systems may be tested, while preventing the disclosure of sensitive test data to participants. At the same time, the test environment also operates as a “black box,” affording participants confidence that their proprietary systems cannot easily be reverse engineered.

The task to be solved in the Total Recall Track is the following:

Given a simple topic description – like those used for ad-hoc and Web search – identify the documents in a corpus, one at a time, such that, as nearly as possible, all relevant documents are identified before all non-relevant documents. Immediately after each document is identified, its ground-truth relevance or non-relevance is disclosed.

Datasets, topics, and automated relevance assessments were all provided by a Web server supplied by the Track. Participants were required to implement either a fully automated (“automatic”) or semi-automated (“manual”) process to download the datasets and topics, and to submit documents for assessment to the Web server, which rendered a relevance assessment for each submitted document in real time. Thus, participants were tasked with identifying documents for review, while the Web server simulated the role of a human-in-the-loop assessor operating in real time. Rank-based and set-based evaluation measures were calculated based on the order in which documents were presented to the Web server for assessment, as well as the set of documents that were presented to the Web server at the time a participant “called their shot,” or declared that a “reasonable” result had been achieved. Particular emphasis was placed on achieving high recall while reviewing the minimum possible number of documents.

The TREC 2015 Total Recall Track used a total of eight test collections: three for Practice runs, three for “At-Home” participation, and two for “Sandbox” participation. Practice and At-Home participation were done using the open Web: Participants ran their own systems and connected to the Web server at a public address. The Practice collections were available for several weeks prior to the At-Home collections; the At-Home collections were available for official runs throughout July and August 2015 (and continue to be available for unofficial runs).

Sandbox runs were conducted entirely on a Web-isolated platform hosting the data collections, from mid-September through mid-November 2015. To participate in the Sandbox task, participants were required to encapsulate – as a VirtualBox virtual machine – a fully autonomous solution that would contact the Web server and conduct the task without human intervention. The only feedback available to Sandbox participants consisted of summary evaluation measures showing the number of relevant documents identified, as a function of the total number of documents identified to the Web server for review.

To aid participants in the Practice, At-Home, and Sandbox tasks, as well as to provide a baseline for comparison, a Baseline Model Implementation (“BMI”) was made available to participants.¹ BMI was run on all of the

*Current affiliation: University of Waterloo. The views expressed herein are solely those of the author and should not be attributed to her former firm or its clients.

¹<http://plg.uwaterloo.ca/~gvcormac/trecvm/>.

collections, and summary results were supplied to participants for their own runs, as well as for the BMI runs. The system architecture for the Track is detailed in a separate Notebook paper titled *Total Recall Track Tools Architecture Overview* [16].

The TREC 2015 Total Recall Track attracted 10 participants, including three industrial groups that submitted “manual athome” runs, two academic groups that submitted only “automatic athome” runs, and five academic groups that submitted both “automatic athome” and “sandbox” runs.

The 2015 At-Home collections consisted of three datasets and 30 topics. The Jeb Bush emails² were collected and assessed for 10 topics by the Track coordinators. The “Illicit Goods” and “Local Politics” datasets, along with 10 topics for each, were derived from the Dynamic Domain datasets³ and assessed by the Total Recall coordinators. These collections continue to be available through the Total Recall Server to 2015 participants, and were made available to 2016 participants for training purposes.

The Sandbox collections consisted of two datasets and 23 topics. On-site access to former Governor Tim Kaine’s email collection at the Library of Virginia⁴ was arranged by the Track coordinators, where a “Sandbox appliance” was used to conduct and evaluate participant runs according to topics that corresponded to archival category labels previously applied by the Library’s Senior State Records Archivist: “Not a Public Record,” “Open Public Record,” “Restricted Public Record,” and “Virginia Tech Shooting Record.” The coordinators also secured approval to use the MIMIC II clinical dataset⁵ as the second Sandbox dataset. The textual documents from this dataset – consisting of discharge summaries, nurses’ notes, and radiology reports – were used as the corpus; the 19 top-level codes in the ICD-9 hierarchy⁶ were used as the “topics.”

The principal tool for comparing runs was a *gain curve*, which plots recall (*i.e.*, the proportion of all relevant documents submitted to the Web server for review) as a function of effort (*i.e.*, the total number of documents submitted to the Web server for review). A run that achieves higher recall with less effort demonstrates superior effectiveness, particularly at high recall levels. The traditional *recall-precision* curve conveys similar information, plotting precision (*i.e.*, the proportion of documents submitted to the Web server that are relevant) as a function of recall (*i.e.*, the proportion of all relevant documents submitted to the Web server for review). Both curves convey similar information, but are influenced differently by prevalence or richness (*i.e.*, the proportion of documents in the collection that are relevant), and convey different impressions when averaged over topics with different richness.

A gain curve or recall-precision curve is blind to the important consideration of when to stop a retrieval effort. In general, the density of relevant documents diminishes as effort increases, and at some point, the benefit of identifying more relevant documents no longer justifies the review effort required to find them. Participants were asked to “call their shot,” or to indicate when they thought a “reasonable” result had been achieved; that is, to specify the point at which they would recommend terminating the review process because further effort would be “disproportionate.” They were not actually required to stop at this point, they were simply given the option to indicate, contemporaneously, when they would have chosen to stop had they been required to do so. For this point, we report traditional set-based measures such as recall, precision, and F_1 .

To evaluate the appropriateness of various possible stopping points, the Track coordinators devised a new parametric measure: *recall @ $aR + b$* , for various values of a and b . *Recall @ $aR + b$* is defined to be the recall achieved when $aR + b$ documents have been submitted to the Web server, where R is the number of relevant documents in the collection. In its simplest form *recall @ $aR + b$* [$a = 1; b = 0$] is equivalent to R-precision, which has been used since TREC 1 as an evaluation measure for relevance ranking. R-precision might equally well be called R-recall, as precision and recall are, by definition, equal when R documents have been reviewed. The parameters a and b allow us to explore the recall that might be achieved when a times as many documents, plus an additional b documents are reviewed. The parameter a admits that it may be reasonable to review more than one document for every relevant one that is found; the parameter b admits that it may be reasonable to review a fixed number of additional documents, over and above the number that are relevant. For example, if there are 100 relevant documents in the collection, it may be reasonable to review 200 documents ($a = 2$), plus an additional 100 documents ($b = 100$), for a total of 300 documents, in order to achieve high recall. In this Track Overview paper, we report all combinations of $a \in \{1, 2, 4\}$ and $b \in \{0, 100, 1000\}$.

At the time of 2015 Total Recall Track, the coordinators had hoped to be able to implement *facet*-based variants of the recall measures described above (*see* Cormack & Grossman [3]), but suitable relevance assessments for the facets were not available in time. We therefore decided to implement such measures in a future Track. The rationale for facet-based measures derives from the fact that, due to a number of factors including assessor disagreement,

²<https://web.archive.org/web/20160221072908/http://jebemails.com/home>.

³<http://trec-dd.org/>.

⁴<http://www.virginiamemory.com/collections/kaine/>.

⁵<https://physionet.org/mimic2/>.

⁶https://en.wikipedia.org/wiki/List_of_ICD-9_codes.

a goal of recall=1.0 is neither reasonable nor achievable in most circumstances. However, it is difficult to justify an arbitrary lower target of, say, recall=0.8, without characterizing the nature of the 20% relevant documents that are omitted by such an effort. Are these omitted documents simply marginal documents about whose relevance reasonable people might disagree, or do they represent a unique and important (though perhaps rare) class of clearly relevant documents? To explore this question, we wanted to be able to calculate the recall measures for a given run separately for each of several *facets* representing different classes of documents; a superior high-recall run should be expected to achieve high recall on all facets. This issue remains to be explored.

In calculating effort and precision, the measures outlined above consider only the number of documents submitted to the Web server for assessment. For manual runs, however, participants were permitted to look at the documents, and hence conduct their own assessments. Participants were asked to track and report the number of documents they reviewed; when supplied by participants, these numbers are reported in this Overview and should be considered when comparing manual runs to one another, or to automatic runs. It is not obvious how one would incorporate this effort formulaically into the gain curves, precision-recall curves, and *recall @ aR + b* measures; therefore, the coordinators have chosen not to try.

Results for the TREC 2015 Total Recall Track show that a number of methods achieved results with very high recall and precision, on all collections, according to the standards set by previous TREC tasks. This observation should be interpreted in light of the fact that runs were afforded an unprecedented amount of relevance feedback, allowing them to receive authoritative relevance assessments throughout the process.

Overall, no run consistently achieved higher recall at lower effort than BMI. A number of runs, including manual runs, automatic runs, and the baseline runs, appeared to achieve similar effectiveness – all near the best on every collection – but with no run consistently bettering the rest on every collection. Thus, The 2015 Total Recall Track had no clear “winner.”

2 Test Collections

Each test collection consisted of a corpus of English-language documents, a set of topics, and a complete set of relevance assessments for each topic. For Practice runs, we used three public document corpora for which topics and relevance assessments were available:

- The *20 Newsgroups Dataset*,⁷ consisting of 18,828 documents from each of 20 newsgroups. We used three of the newsgroup subject categories – “space,” “hockey,” and “baseball” – as the three practice topics in the *test* practice collection.
- The *Reuters-21578 Test Collection*,⁸ consisting of 21,578 newswire documents. We used four of the subject categories – “acquisitions,” “Deutsche Mark,” “groundnut,” and “livestock” – as the four practice topics in the *test* practice collection.
- The *Enron Dataset* used by the TREC 2009 Legal Track [10]. We used a version of this dataset captured by the University of Waterloo in the course of its participation in TREC 2009, modified to exclude vacuous documents, resulting in a corpus of 723,537 documents. We used two of the topics from the TREC 2009 Legal Track – “Fantasy Football” and “Prepay Transactions” – as the two practice topics for the *bigtest* practice collection. The relevance assessments were derived from those rendered by the University of Waterloo team and the official TREC assessments, with deference to the official assessments.

For the At-Home runs, we used three new datasets:

- The (redacted) Jeb Bush Emails,⁹ consisting of 290,099 emails from Jeb Bush’s eight-year tenure as Governor of Florida. We used 10 issues associated with his governorship as topics for the *athome1* test collection: “school and preschool funding,” “judicial selection,” “capital punishment,” “manatee protection,” “new medical schools,” “affirmative action,” “Terri Schiavo,” “tort reform,” “Manatee County,” and “Scarlet Letter Law.” Using the continuous active learning (“CAL”) method of Cormack and Mojdeh [5], the Track coordinators assessed documents in the corpus to identify as many of the relevant documents for each topic as reasonably possible.

⁷<http://qwone.com/~jason/20Newsgroups/>.

⁸<http://www.daviddlewis.com/resources/testcollections/reuters21578/>.

⁹<https://web.archive.org/web/20160221072908/http://jebemails.com/home>.

Participating Team	athome1	athome2	athome3	Kaine	MIMIC II
Catres	1M				
CCRI	1A	1A	1A		
eDiscoveryTeam	1M	1M	1M		
NINJA	1M	1M	1M		
TUW	6A	6A	6A	6A	6A
UvA.ILPS	2A	2A	2A	2A	2A
WaterlooClarke	2A	1A	1A	1A	1A
WaterlooCormack	3A	3A	3A	1A	1A
Webis	2A	2A	2A	2A	2A
WHU_IRGroup	1A				

Table 1: Participation in the TREC 2015 Total Recall Track. Table entries indicate the number of runs submitted for each test collection by a particular participating team. “M” indicates manual runs, “A” indicates automatic runs.

- The *Illicit Goods* dataset collected for the TREC 2015 Dynamic Domain Track [17]. We used 465,147 documents collected from Blackhat World¹⁰ and Hack Forum.¹¹ For the *athome2* test collection, we used 10 of the many topics that were composed and partially assessed by NIST assessors for use by the Dynamic Domain Track: “paying for Amazon book reviews,” “CAPTCHA services,” “Facebook accounts,” “surely Bitcoins can be used,” “Paypal accounts,” “using TOR for anonymous browsing,” “rootkits,” “Web scraping,” “article spinner spinning,” and “offshore Web sites.” The Track coordinators re-assessed the documents using the CAL method of Cormack and Mojdeh [5], to identify as many of the relevant documents for each topic as reasonably possible.
- The *Local Politics* dataset collected for the TREC 2015 Dynamic Domain Track [17]. We used 902,434 articles collected from news sources in the northwestern United States and southwestern Canada. For the *athome3* test collection, we used 10 of the many topics that were composed and partially assessed by NIST assessors for use by the Dynamic Domain Track: “Pickton murders,” “Pacific Gateway,” “traffic enforcement cameras,” “rooster chicken turkey nuisance,” “Occupy Vancouver,” “Rob McKenna gubernatorial candidate,” “Rob Ford Cut the Waist,” “Kingston Mills lock murder,” “fracking,” and “Paul and Cathy Lee Martin.” The Track coordinators re-assessed the documents using the CAL method of Cormack and Mojdeh [5], to identify as many of the relevant documents for each topic as reasonably possible.

For the Sandbox runs, we used two new datasets:

- The Kaine Email Collection at the Library of Virginia.¹² From the 1.3M email messages from Tim Kaine’s eight-year tenure as Governor of Virginia, we used 401,953 that had previously been labeled by the Virginia Senior State Records Archivist according to the following four categories: “public record,” “open record,” “restricted record,” and “Virginia Tech shooting ([subject to a legal] hold).” Each of the four categories was used as a topic in the *kaine* test collection. The runs themselves were executed on an isolated computer installed at the Library of Virginia and operated by Library of Virginia staff.
- The MIMIC II Clinical Dataset,¹³ consisting of anonymized, time-shifted, records for 31,538 patient visits to an Intensive Care Unit. We used the textual record for each patient – consisting of one or more nurses’ notes, radiology reports, and discharge summaries – as a “document” in the corpus, and each of 19 top-level ICD-9 codes supplied with the dataset as a topic for the *mimic* test collection: “infectious and parasitic diseases,” “neoplasms,” “endocrine, nutritional and metabolic diseases, and immunity disorders,” “diseases of the blood and blood-forming organs,” “mental disorders,” “diseases of the nervous system and sense organs,” “diseases of the circulatory system,” “diseases of the respiratory system,” “diseases of the digestive system,” “diseases of the genitourinary system,” “complications of pregnancy, childbirth, and the puerperium,” “diseases of the skin and subcutaneous tissue,” “diseases of the musculoskeletal system and connective tissue,” “congenital anomalies,” “certain conditions originating in the perinatal period,” “symptoms, signs, and ill-defined conditions,” “injury and poisoning,” “factors influencing health status and contact with health services,” and

¹⁰<http://www.blackhatworld.com/>.

¹¹<http://hackforums.net/>.

¹²<http://www.virginiamemory.com/collections/kaine/under-the-hood>.

¹³<https://physionet.org/mimic2/>.

“external causes of injury and poisoning.” The runs were executed on an isolated computer installed at the University of Waterloo and operated by the Track coordinators.

Table 1 shows the number of runs submitted for each test collection by each participating team.

3 Participant Submissions

The following descriptions are paraphrased from responses to a required questionnaire submitted by each participating team.

3.1 UvA.ILPS

The UvA.ILPS team [6] used automatic methods for the At-Home and Sandbox tests that modified the Baseline Model Implementation in two ways:

1. adjusted the batch size based on the number of retrieved relevant documents and stopped after a threshold if the batch contained no relevant documents;
2. one run used logistic regression (as per the Baseline Model Implementation), while another used random forests.

3.2 WaterlooClarke

The WaterlooClarke team [9] used automatic methods for At-Home and Sandbox tests that:

1. employed clustering to improve the diversity of feedback to the learning algorithm;
2. used n-gram features beyond the bag-of-words tf-idf model provided in the Baseline Model Implementation;
3. employed query expansion;
4. used the fusion of differently ranking algorithms.

The WaterlooClarke team consisted of a group of graduate students who had no access to the test collections beyond that afforded to all participating teams.

3.3 WaterlooCormack

The WaterlooCormack team [4] employed the Baseline Model Implementation, without modification, except to “call its shot” to determine when to stop. Therefore, the gain curves, recall-precision curves, and *recall @ aR + b* statistics labeled “WaterlooCormack” are synonymous with the Baseline Model Implementation.

Two different stopping criteria were investigated:

1. a “knee-detection” algorithm was applied to the gain curve, and the decision that a reasonable result had been achieved was made when the slope of the curve after the knee was a fraction of the slope before the knee;
2. a “reasonable” result was deemed to have been achieved when m relevant and n non-relevant documents had been reviewed, where $n > a \cdot m + b$, where a and b are predetermined constants. For example, when $a = 1$ and $b = 2399$, review would be deemed to be complete when the number of non-relevant documents retrieved was equal to the number of relevant documents retrieved, plus 2,399. In general, the constant a determines how many non-relevant documents are to be reviewed in the course of finding each relevant document, while b represents fixed overhead, independent of the number of relevant documents.

The WaterlooCormack team consisted of Gordon V. Cormack and Maura R. Grossman, who were both Track coordinators. The Baseline Model Implementation was fixed prior to the development of any of the datasets. Cormack and Grossman had knowledge of the At-Home test collections, but not the Sandbox test collections, when the stopping criteria were chosen.

3.4 Webis

The Webis team [15] employed two methods:

1. a basic naïve approach in retrieving as many relevant documents as possible;
2. a keyphrase experiment that built on the BMI system by intelligently obtaining a list of phrases from documents judged by the API as relevant and using them as new topics for ad-hoc search.

3.5 CCRi

The CCRi team [7] represented words from the input corpus as vectors using a neural network model and represented documents as a tf-idf weighted sum of their word vectors. This model was designed to produce a compact versions of tf-idf vectors while incorporating information about synonyms. For each query topic, CCRi attached a neural network classifier to the output of BMI. Each classifier was updated dynamically with respect to the given relevance assessments.

3.6 eDiscoveryTeam

eDiscoveryTeam [13] employed a manual approach for the athome1, athome2, and athome3 test collections. Eight hours of manual search and review were conducted, on average, per topic. Two of the eight hours were spent composing 25 queries (per topic, on average) and examining their results; six hours were spent reviewing 500 documents (per topic, on average), of which *only those deemed relevant were submitted to the automated assessment server*. During the search and review process, Web searches were conducted where necessary to inform the searchers.

3.7 NINJA

The NINJA team employed a manual approach for the athome1, athome2, and athome3 test collections. One hour of manual search was conducted, in which three queries were composed, on average, per topic. Wikipedia and Google searches were used to inform the searchers. A commercial “predictive coding” tool, trained using the results of the queries, was used to generate the NINJA runs. No documents from the test collection were reviewed prior to being submitted to the automated assessment server.

3.8 catres

The catres team [12] employed a manual approach for the athome1 test collection. A group of searchers independently spent one hour each investigating each topic, after which a learning tool was used to generate the run. An average of eight manual queries were used per topic, and an average of 262 documents were reviewed. Every document reviewed by the team was also submitted to the automated assessment server.

3.9 TUW

The TUW team [14] applied six variants on the Baseline Model Implementation to all of the At-Home and Sandbox test collections. The variants included the use and non-use of a BM25 ranking, the use and non-use of stop words, and the use and non-use of tf-idf weighting.

3.10 WHU_IRGroup

The WHU_IRGroup team [1] applied iterative query expansion to the athome1 test collection.

4 Results

4.1 Gain Curves and Recall-Precision Curves

Figure 1 plots the effectiveness of the best run for each of the participating teams on the athome1 test collection. The top panel plots effectiveness as a gain curve, while the bottom panel plots effectiveness as a recall-precision curve. The gain curve shows that a number of the systems achieved 90% recall, on average, with a review effort of 10,000 documents. The recall-precision curve, on the other hand, shows that a number of the systems achieved 90%

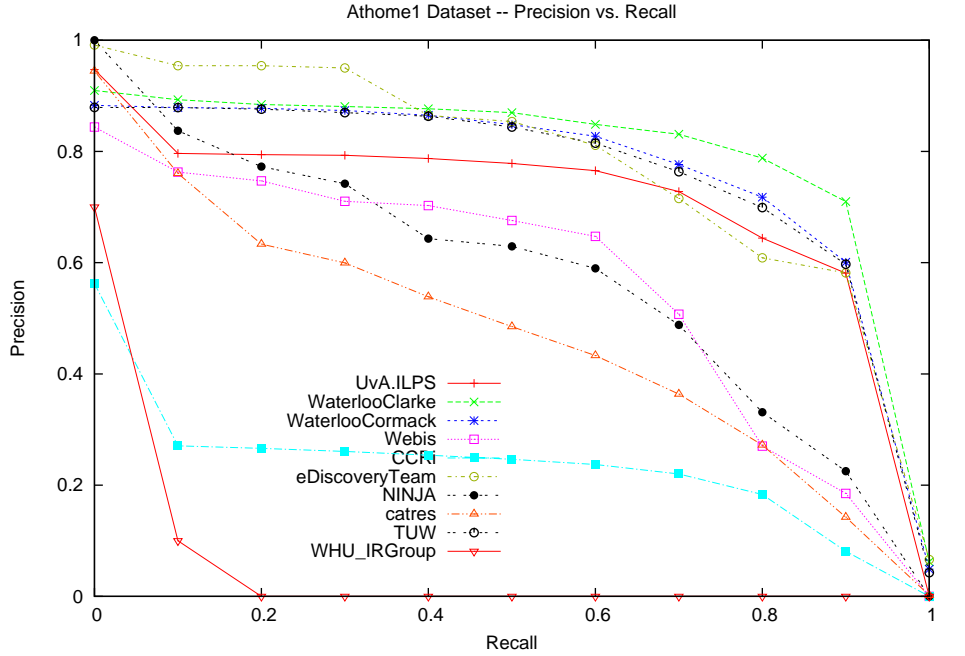
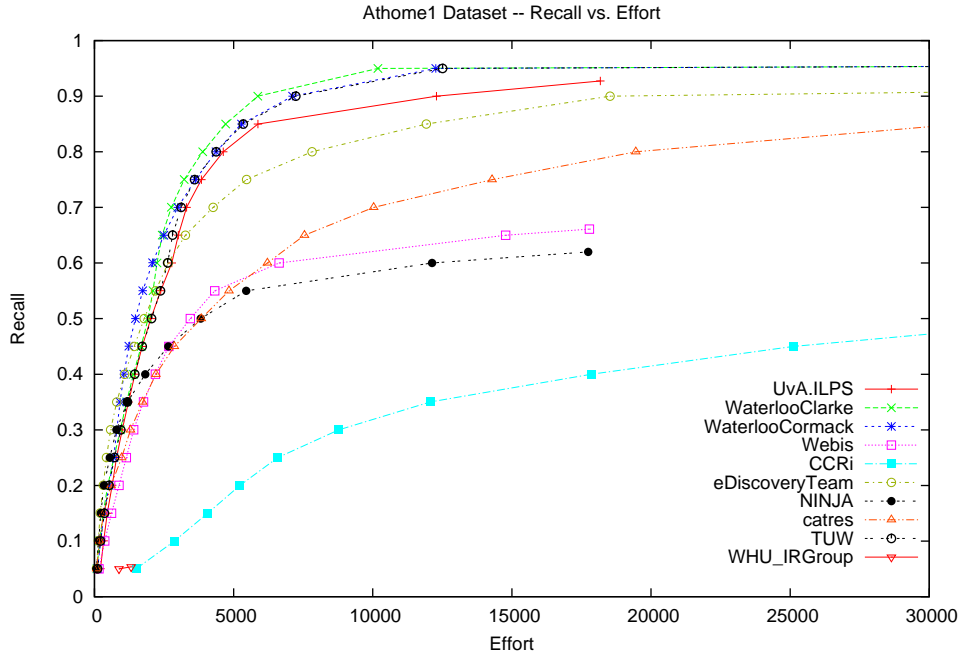


Figure 1: Athome1 Results – Average Gain and Interpolated Recall-Precision Curves.

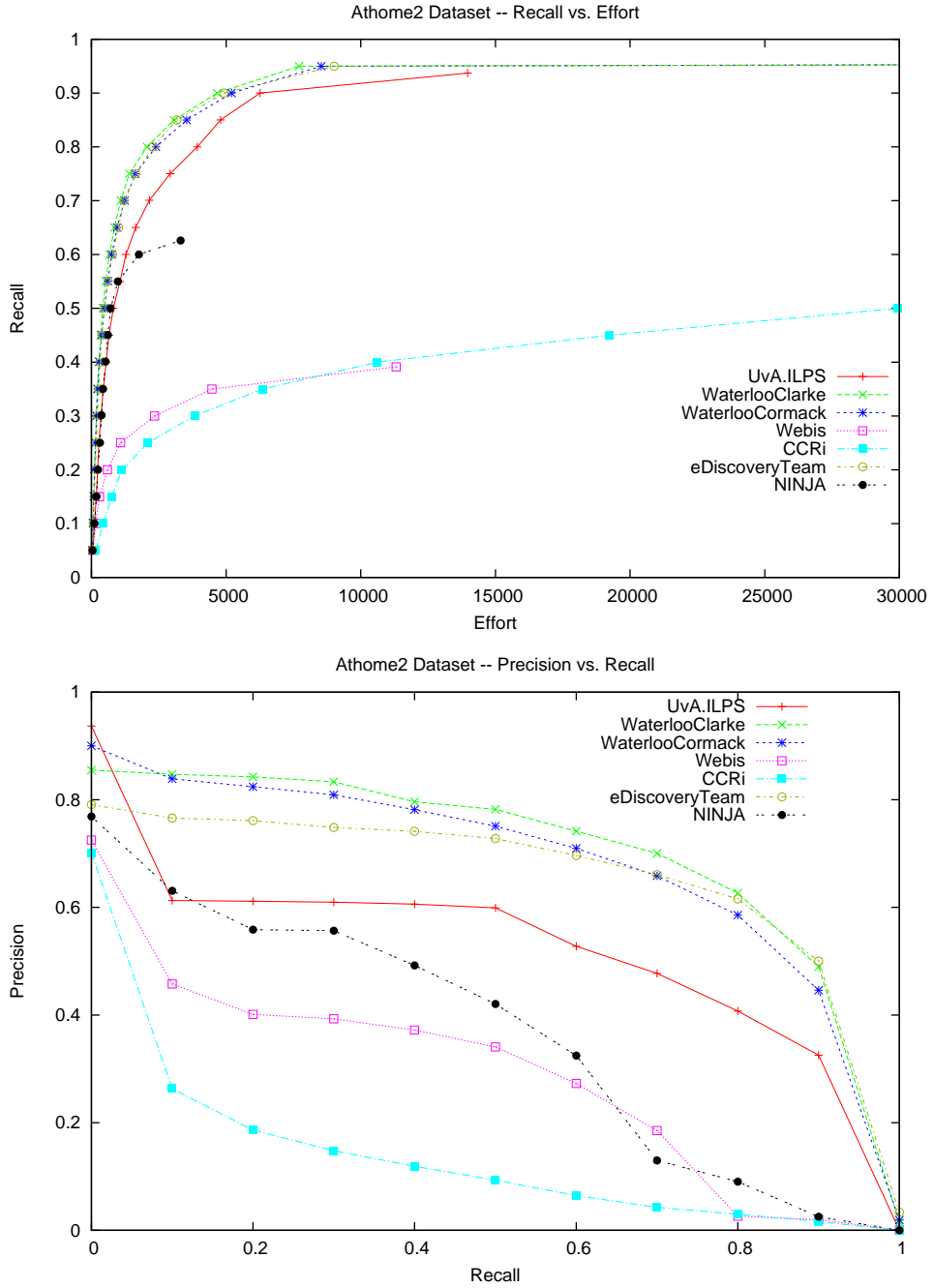


Figure 2: Athome2 Results – Average Gain and Interpolated Recall-Precision Curves.

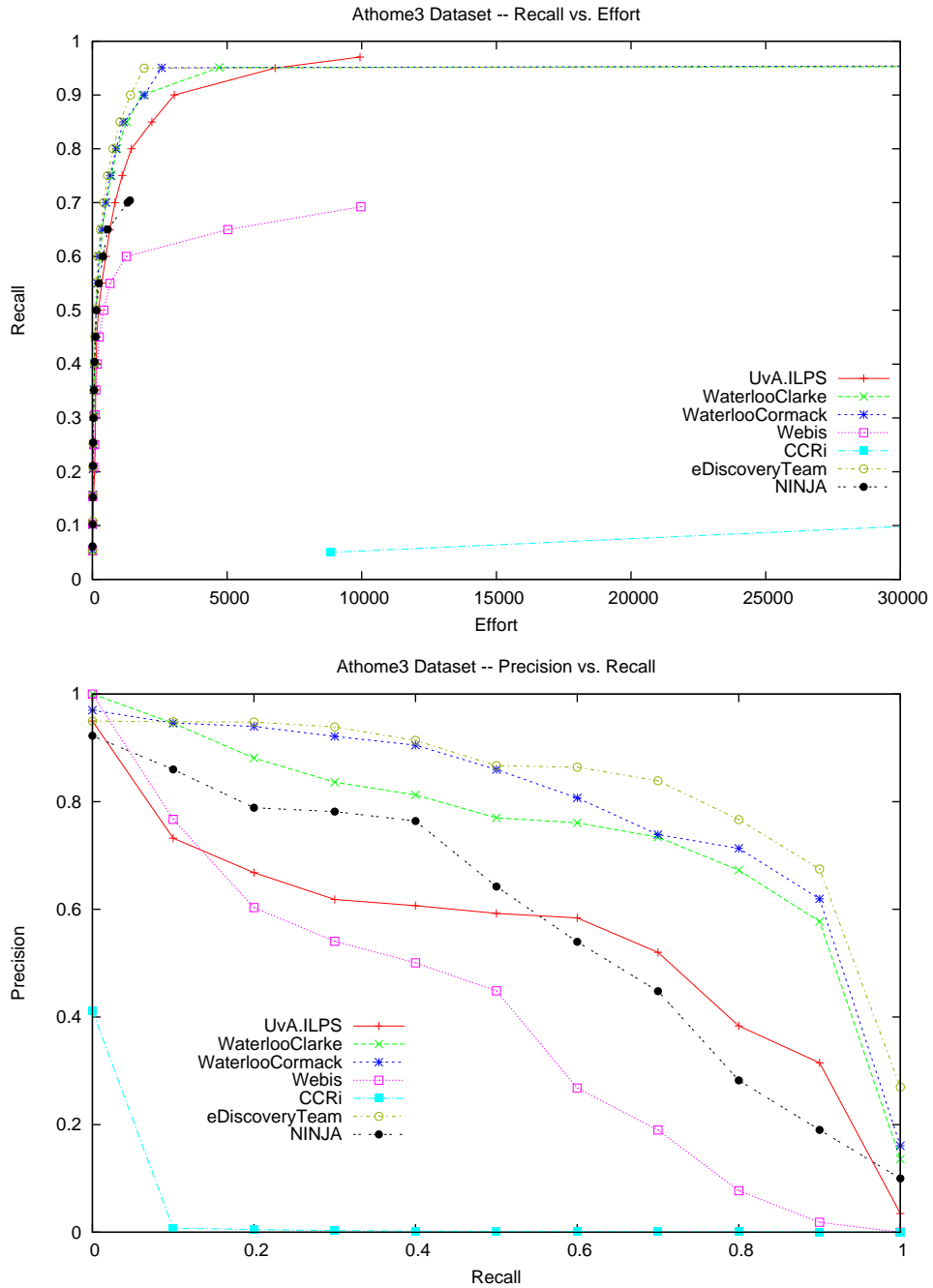


Figure 3: Athome3 Results – Average Gain and Interpolated Recall-Precision Curves.

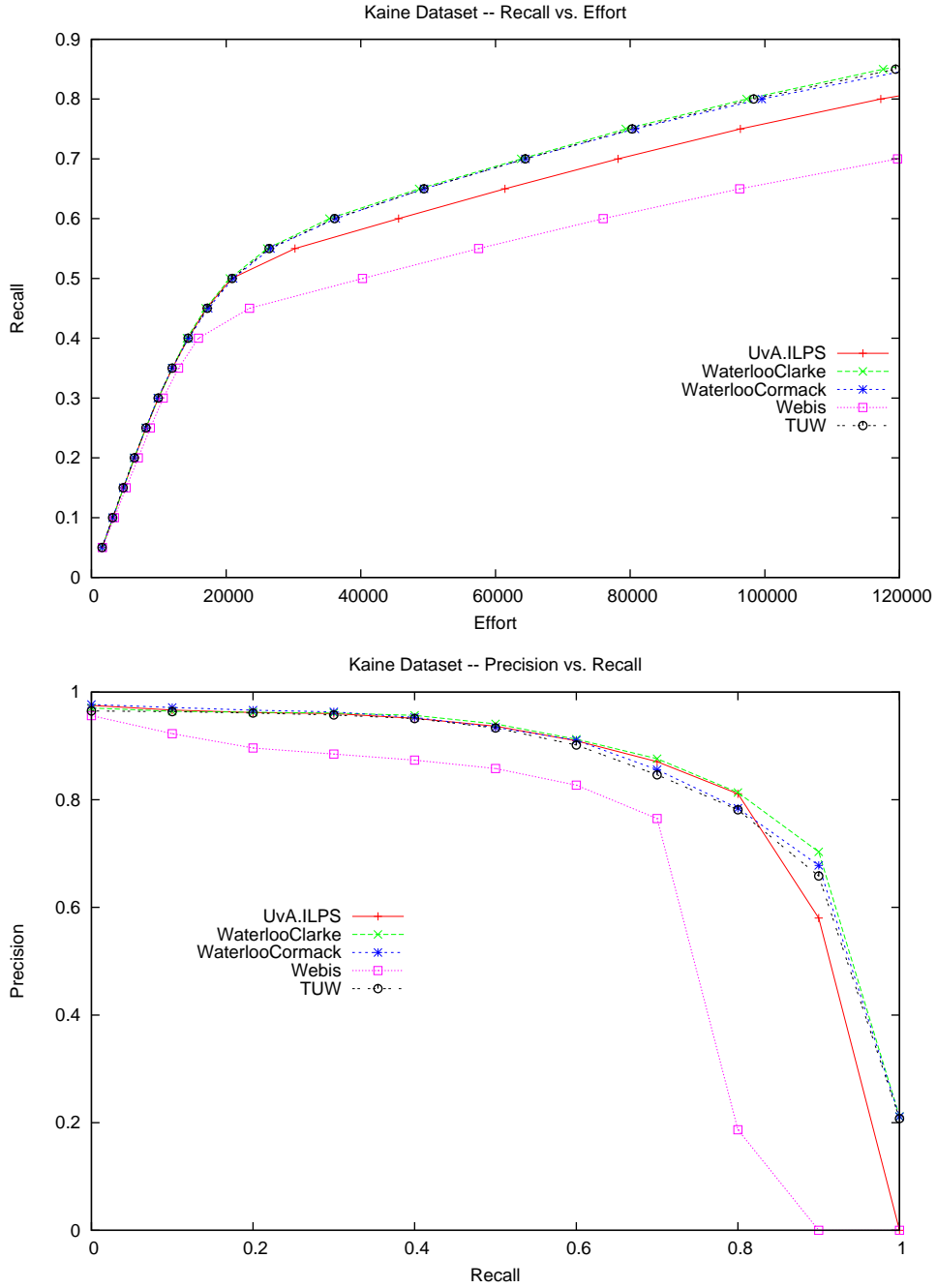


Figure 4: Kaine Results – Average Gain and Interpolated Recall-Precision Curves.

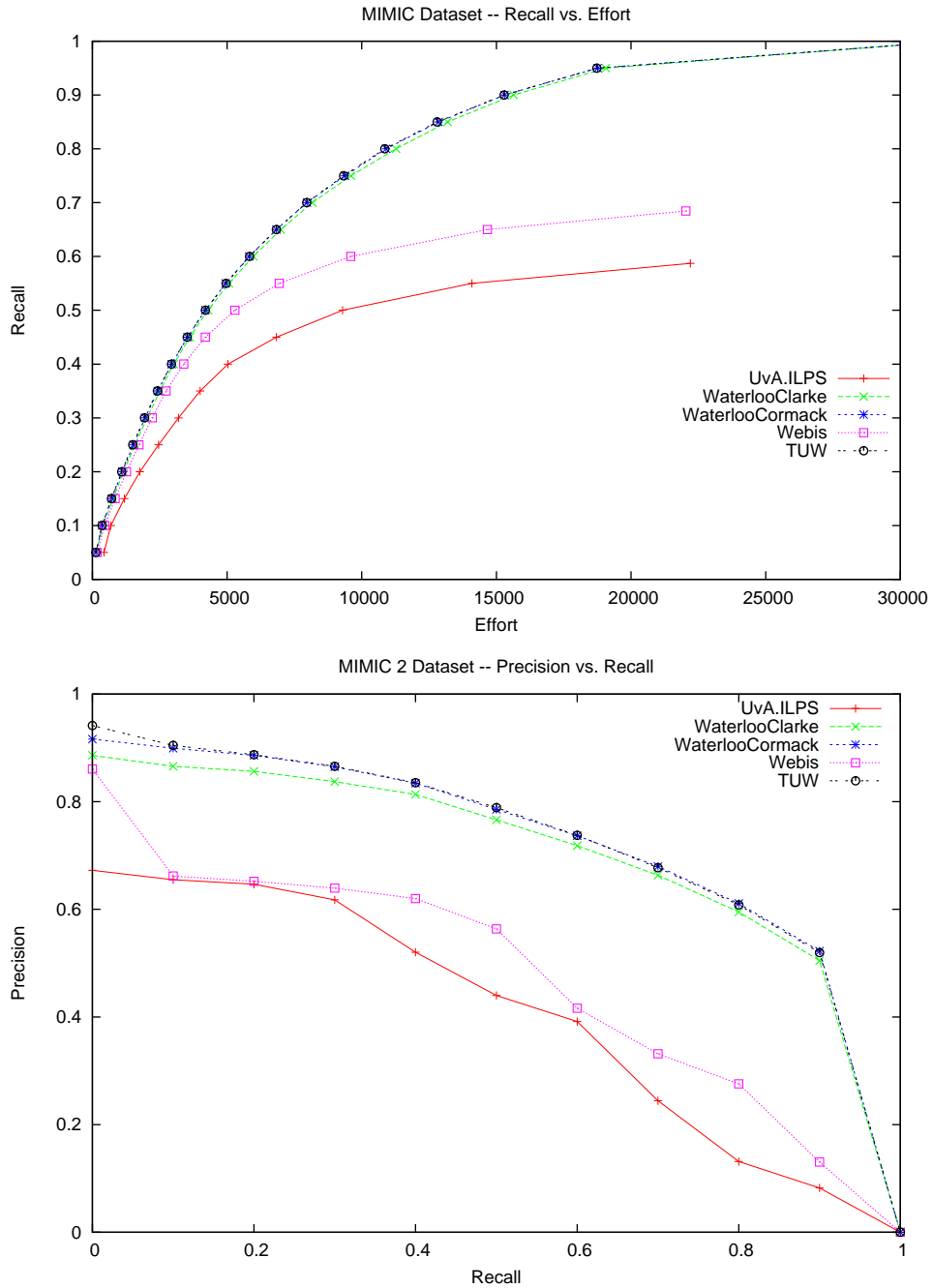


Figure 5: MIMIC II Results – Average Gain and Interpolated Recall-Precision Curves.

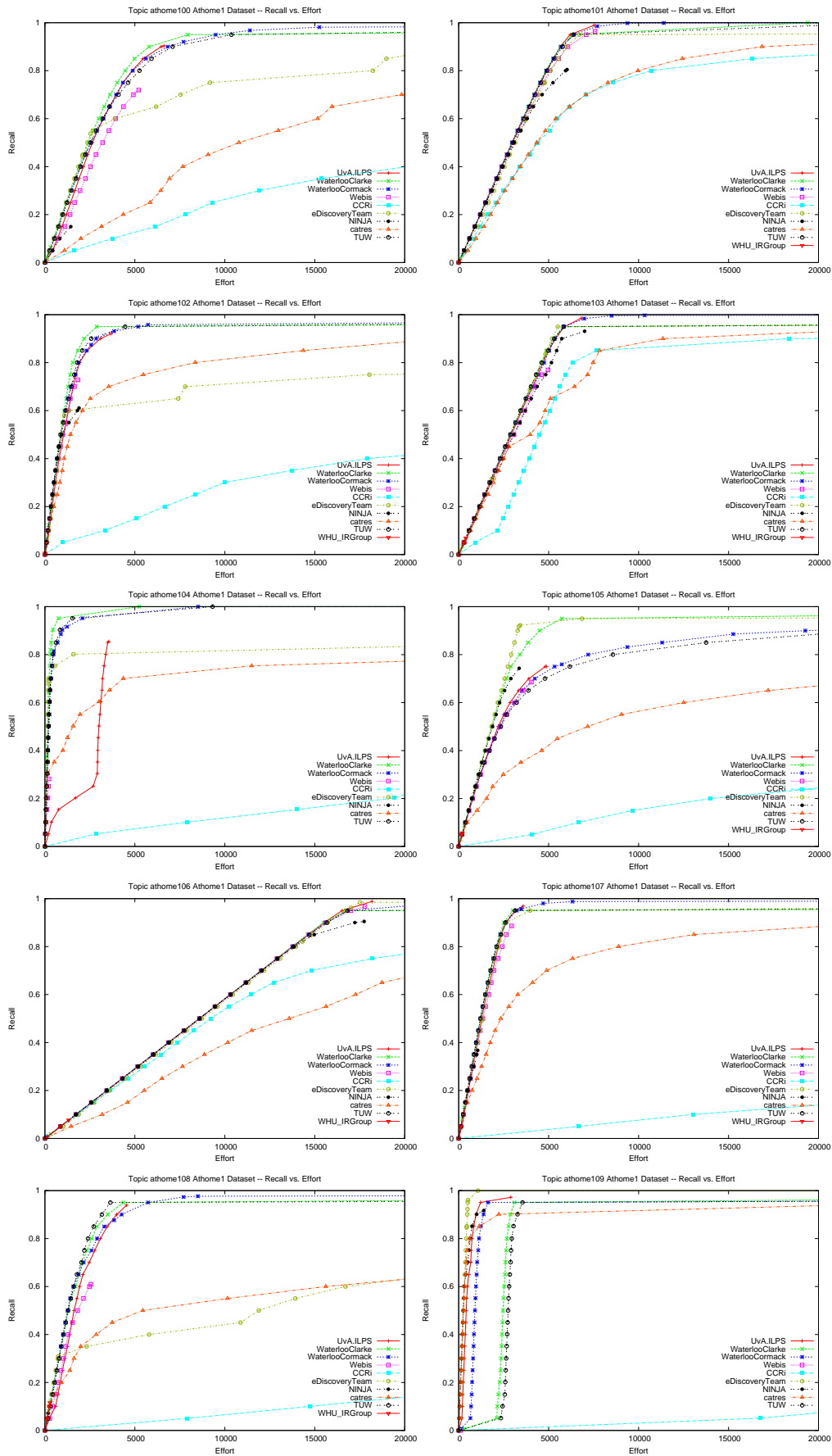


Figure 6: Per-Topic Gain Curves for the Athome1 Test Collection.

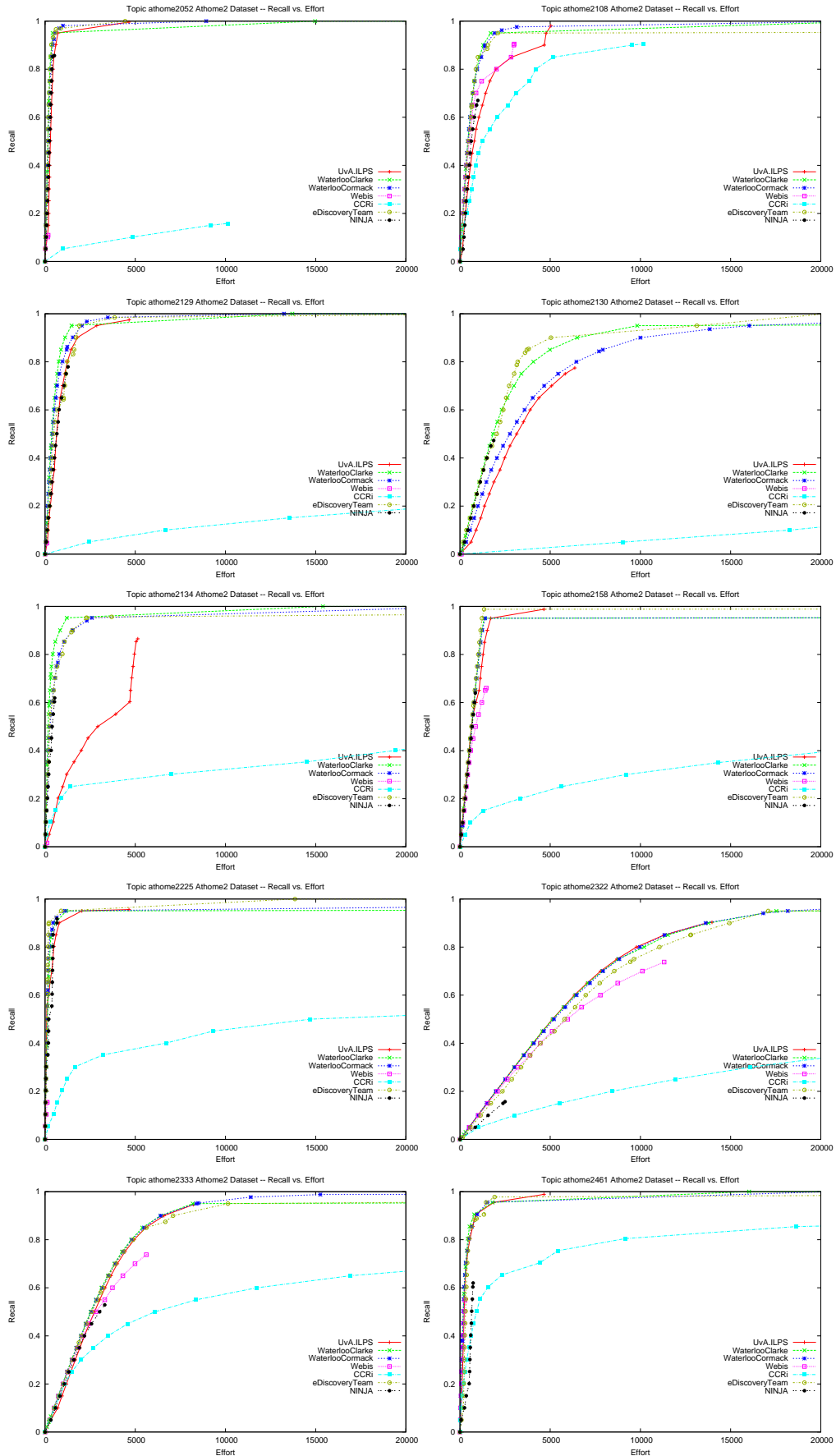


Figure 7: Per-Topic Gain Curves for the Athome2 Test Collection.

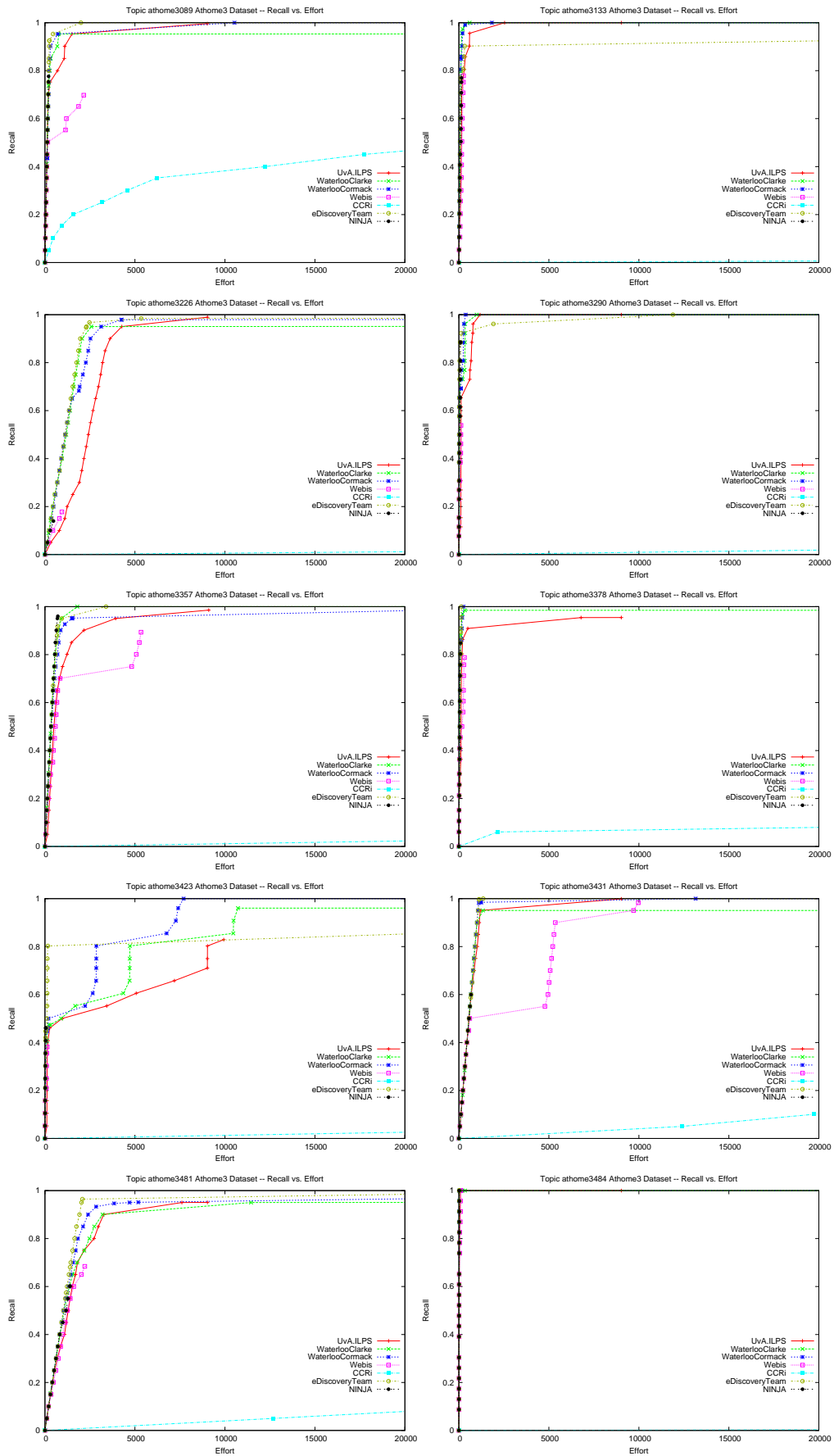


Figure 8: Per-Topic Gain Curves for the Athome3 iTest Collection.

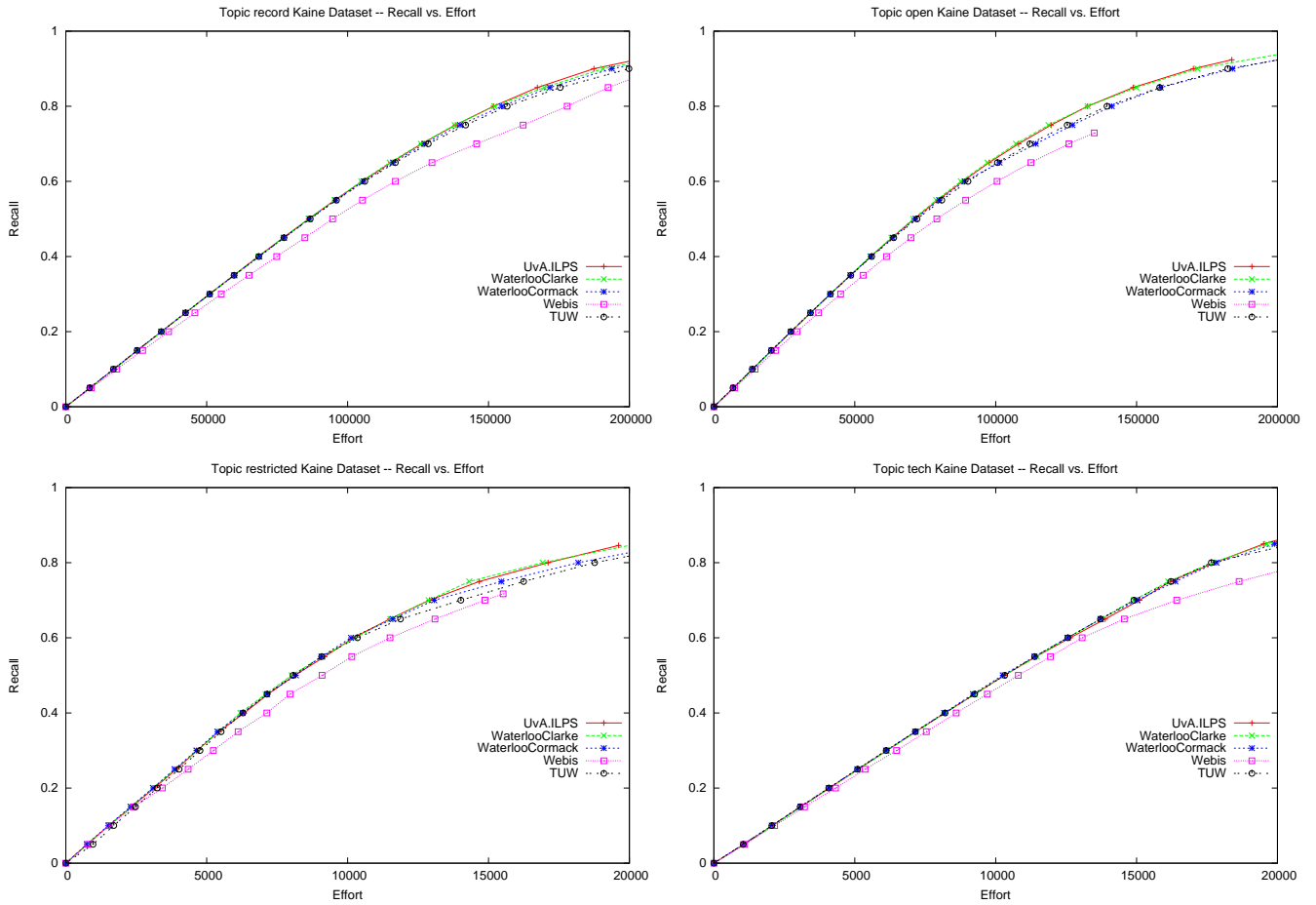


Figure 9: Per-Topic Gain Curves for the Kaine Test Collection.

recall, with higher than 50% precision, on average. It is not obvious which of these results is preferable. For a topic with only 100 relevant documents, it may be considered unreasonable to require the review of 10,000 documents; the average gain curve gives little insight in this regard. Individual per-topic gain curves (Figure 6) offer per-topic insight, but are challenging to generalize. The recall-precision curve indicates that it is possible to achieve 90% recall, with 50% precision or better; but if there are 10,000 documents to be found, is 50% good enough? Moreover, is a system that requires the review of 1,000 documents to find 100 (*i.e.*, 10% precision) inferior to a system that requires the review of 5,000 documents to find 1,000 (*i.e.*, 20% precision)? We suggest that a reasonable measure may lie somewhere in between these two extremes: The effectiveness of a system depends both on the absolute effort required (as shown by effort in the gain curve) and the effort required relative to the number of relevant documents found (as shown by precision in the recall-precision curve). Ideally, an effective system would score well on both curves.

The average gain and recall-precision curves for the remaining Athome test collections – athome2 and athome3 – are shown in Figures 2 and 3; per-topic gain curves are shown in Figures 7 and 8. Average curves for the Sandbox test collections – Kaine and MIMIC II – are shown in Figures 4 and 5. Per-topic gain curves for the Kaine collection are shown in Figure 9; per-topic curves for MIMIC-II are omitted for brevity. From these curves it can be seen that a number of systems – including the Baseline Model Implementation (denoted as “WaterlooCormack” in the curves) – achieve similar effectiveness, but no single system dominates over all topics, or over all collections. Further study is necessary to determine whether the observed differences among the top-performing systems on particular topics and test collections represent real, reproducible differences.

4.2 Recall @ $aR+b$

Tables 2 through 6 show, for each test collection, the new measure *Recall @ $aR+b$* . This measure quantifies the tradeoff between achieving high recall with effort proportionate to the number of relevant documents, and achieving high recall with reasonable overall effort.

4.3 When to Stop?

Participants were afforded the opportunity to “call their shot,” indicating the point at which they would have recommended terminating the search because the additional effort to identify more relevant documents would have been unreasonable or disproportionate. Some participants did “call their shot,” while others simply terminated their retrieval effort. In both cases, we tabulated the recall that had been achieved at that point and the effort (in terms of the number of assessed documents) necessary to achieve it. Clearly, there is a tradeoff between recall and effort; we made no attempt to quantify what was a reasonable compromise, and instead, present the raw, per-topic precision and effort results in Figures 7 through 12.

5 Discussion

In 2015, the inaugural year of the Total Recall Track, it was necessary to develop new, completely labeled datasets, a new evaluation architecture, new evaluation measures, and a Baseline Model Implementation.

For the Athome task, the original plan was to use datasets labeled by NIST assessors on a five-point relevance scale, according to a two-level hierarchy of topics and subtopics. This labeling effort was not completed within the allotted timeframe and budget, so the Total Recall coordinators decided to conduct binary relevance assessments for 30 topics across three datasets. Candidate documents for assessment were selected using a combination of ad-hoc search and machine learning. With few exceptions (*e.g.*, documents containing “schivo” or “lethal injection”), every document labeled relevant was assessed by one of the Track coordinators. The selection and assessment of documents continued until the coordinators believed that substantially all relevant documents had been found; documents that were not selected for assessment were summarily labeled “not relevant.” The Sandbox task used binary relevance assessments that did not rely in any way on the Track coordinators, or any search or machine-learning software. The Kaine collection was exhaustively labeled by the Senior State Records Archivist; the MIMIC II collection contained ICD-9 diagnostic codes which were used as relevance labels. The similarity between the Athome and Sandbox results suggests that any confounding due to the method of selecting and assessing the documents was minimal.

Run	Recall @									
	R	R+100	R+1000	2R	2R+100	2R+1000	4R	4R+100	4R+1000	4R+1000
catres	0.496	0.505	0.569	0.654	0.657	0.691	0.746	0.748	0.771	0.771
CCRi	0.233	0.236	0.260	0.309	0.310	0.320	0.364	0.364	0.372	0.372
eDiscoveryTeam	0.783	0.792	0.819	0.842	0.843	0.853	0.862	0.862	0.873	0.873
NINJA	0.577	0.600	0.620	0.619	0.619	0.620	0.620	0.620	0.620	0.620
TUW-INB	0.710	0.740	0.815	0.841	0.843	0.861	0.870	0.871	0.955	0.955
TUW-1SB	0.714	0.737	0.814	0.838	0.842	0.859	0.869	0.873	0.966	0.966
TUW-1ST	0.715	0.735	0.898	0.886	0.905	0.952	0.966	0.966	0.972	0.972
TUW-6NB	0.717	0.741	0.817	0.841	0.844	0.860	0.871	0.872	0.961	0.961
TUW-6SB	0.630	0.654	0.734	0.753	0.762	0.778	0.789	0.791	0.874	0.874
TUW-6ST	0.710	0.727	0.892	0.866	0.888	0.950	0.964	0.965	0.972	0.972
UvA_ILPS-baseline	0.721	0.740	0.837	0.846	0.850	0.861	0.858	0.859	0.863	0.863
UvA_ILPS-baseline2	0.452	0.467	0.595	0.669	0.678	0.722	0.727	0.727	0.735	0.735
WaterlooClarke-UWPAH1	0.762	0.784	0.850	0.876	0.879	0.887	0.892	0.897	0.988	0.988
WaterlooClarke-UWPAH2	0.761	0.781	0.848	0.868	0.874	0.900	0.910	0.922	0.987	0.987
WaterlooCormack-stop2399	0.716	0.739	0.901	0.903	0.916	0.955	0.968	0.968	0.973	0.973
WaterlooCormack-Knee100	0.707	0.736	0.904	0.901	0.917	0.956	0.968	0.970	0.974	0.974
WaterlooCormack-Knee1000	0.715	0.738	0.900	0.903	0.917	0.953	0.967	0.968	0.973	0.973
webis-baseline	0.628	0.635	0.659	0.661	0.661	0.661	0.661	0.661	0.661	0.661
webis-keyphrase	0.560	0.568	0.588	0.590	0.590	0.590	0.590	0.590	0.590	0.590
WHU_IRGroup	0.038	0.038	0.038	0.038	0.038	0.038	0.038	0.038	0.038	0.038

Table 2: Recall @ aR+b for the athome1 Test Collection.

Run	Recall @								
	R	R+100	R+1000	2R	2R+100	2R+1000	4R	4R+100	4R+1000
CCRI	0.151	0.167	0.261	0.216	0.230	0.300	0.302	0.308	0.349
eDiscoveryTeam	0.647	0.732	0.893	0.859	0.891	0.951	0.942	0.948	0.969
NINJA	0.414	0.472	0.626	0.535	0.576	0.626	0.626	0.626	0.626
UvA_ILPS-baseline	0.501	0.588	0.799	0.735	0.765	0.854	0.837	0.844	0.873
UvA_ILPS-baseline2	0.260	0.279	0.468	0.429	0.450	0.617	0.615	0.630	0.733
WaterlooClarke-UWPAH1	0.700	0.773	0.915	0.896	0.915	0.959	0.960	0.964	0.978
WaterlooCormack-Knee100	0.670	0.735	0.883	0.856	0.877	0.939	0.941	0.947	0.971
WaterlooCormack-Knee1000	0.675	0.740	0.884	0.854	0.875	0.940	0.943	0.947	0.971
WaterlooCormack-stop2399	0.680	0.744	0.884	0.858	0.875	0.939	0.941	0.945	0.971
webis-baseline	0.328	0.356	0.377	0.377	0.377	0.381	0.382	0.383	0.391
webis-keyphrase	0.323	0.332	0.353	0.352	0.353	0.358	0.358	0.359	0.364

Table 3: Recall @ aR+b for the athome2 Test Collection.

Run	Recall @									
	R	R+100	R+1000	2R	2R+100	2R+1000	4R	4R+100	4R+1000	4R+1000
CCRI	0.009	0.012	0.023	0.016	0.017	0.026	0.024	0.024	0.031	0.031
eDiscoveryTeam	0.823	0.923	0.962	0.934	0.957	0.963	0.958	0.963	0.966	0.966
NINJA	0.602	0.704	0.704	0.670	0.704	0.704	0.704	0.704	0.704	0.704
UvA_ILPS-baseline	0.469	0.720	0.880	0.721	0.833	0.919	0.802	0.872	0.926	0.926
UvA_ILPS-baseline2	0.270	0.324	0.522	0.394	0.429	0.640	0.550	0.567	0.714	0.714
WaterlooClarke-UWPAH1	0.736	0.851	0.927	0.873	0.892	0.932	0.902	0.902	0.934	0.934
WaterlooCormack-Knee100	0.776	0.843	0.935	0.877	0.904	0.942	0.914	0.922	0.948	0.948
WaterlooCormack-Knee1000	0.793	0.857	0.936	0.889	0.909	0.942	0.914	0.922	0.948	0.948
WaterlooCormack-stop2399	0.786	0.859	0.935	0.886	0.913	0.943	0.914	0.922	0.948	0.948
webis-baseline	0.432	0.558	0.619	0.513	0.587	0.620	0.571	0.613	0.670	0.670
webis-keyphrase	0.412	0.444	0.456	0.435	0.452	0.458	0.457	0.459	0.513	0.513

Table 4: Recall @ aR+b for the athome3 test collection.

Run	Recall @									
	R	R+100	R+1000	2R	2R+100	2R+1000	4R	4R+100	4R+1000	
TUW-1NB	0.787	0.788	0.799	0.962	0.963	0.965	0.990	0.990	0.990	
TUW-6SB	0.798	0.799	0.808	0.967	0.967	0.968	0.991	0.991	0.991	
TUW-6ST	0.803	0.805	0.816	0.969	0.969	0.970	0.992	0.992	0.992	
UvA,ILPS-baseline	0.812	0.813	0.824	0.921	0.921	0.921	0.921	0.921	0.921	
WaterlooClarke-UWPAH	0.812	0.813	0.824	0.972	0.972	0.973	0.993	0.993	0.993	
WaterlooCormack-stop2399	0.795	0.797	0.808	0.968	0.969	0.969	0.991	0.992	0.992	
Webis-baseline	0.736	0.738	0.748	0.781	0.781	0.781	0.781	0.781	0.781	
Webis-keyphrase	0.556	0.558	0.568	0.603	0.603	0.603	0.603	0.603	0.603	

Table 5: Recall @ aR+b for the Kaine Test Collection.

Run	Recall @									
	R	R+100	R+1000	2R	2R+100	2R+1000	4R	4R+100	4R+1000	
TUW-1NB	0.688	0.702	0.748	0.894	0.896	0.912	0.973	0.974	0.978	
TUW-1SB	0.691	0.704	0.749	0.894	0.896	0.911	0.973	0.973	0.977	
TUW-1ST	0.695	0.706	0.752	0.896	0.899	0.913	0.974	0.974	0.977	
TUW-6NB	0.681	0.695	0.741	0.889	0.891	0.907	0.971	0.971	0.975	
TUW-6SB	0.685	0.697	0.743	0.890	0.891	0.906	0.971	0.971	0.975	
TUW-6ST	0.689	0.700	0.744	0.891	0.893	0.907	0.972	0.973	0.976	
UvA.IILPS-baseline	0.453	0.478	0.489	0.490	0.490	0.490	0.490	0.490	0.490	
UvA.IILPS-baseline2	0.516	0.518	0.570	0.554	0.565	0.587	0.587	0.587	0.587	
WaterlooClarke-UWPAH1	0.684	0.695	0.742	0.889	0.892	0.908	0.970	0.971	0.975	
WaterlooCormack-stop2399	0.697	0.707	0.753	0.896	0.898	0.913	0.973	0.974	0.977	
webis-baseline	0.624	0.637	0.662	0.680	0.681	0.684	0.684	0.684	0.684	
webis-keyphrase	0.513	0.526	0.539	0.559	0.559	0.559	0.559	0.559	0.559	

Table 6: Recall @ aR+b for the MIMIC II Test Collection.

		Topic (R) – Athome1 Collection									
		100	101	102	103	104	105	106	107	108	109
		(4542)	(5836)	(1624)	(5725)	(227)	(3635)	(17135)	(2375)	(2375)	(506)
eDiscoveryTeam		0.5374	0.7493	0.5800	0.8351	0.6916	0.9224	0.9847	0.8236	0.3091	0.9605
		2536	4771	1071	4817	199	3418	17516	2259	746	510
		+651*	+6841*	+1493*	+7203*	+1091*	+674*	+2226*	+1164*	+696*	+753*
		0.1497	0.8045	0.6108	0.9310	0.7005	0.7431	0.9048	0.3676	0.0720	0.9170
	NINJA	1437	6025	1895	6998	327	3359	17743	1056	171	1396
		0.9036	0.9896	0.9206	0.9857	0.8546	0.7508	0.9879	0.9689	0.9390	0.9723
	UvA_ILPS-baseline	6618	7538	3679	6823	3551	4876	18180	3582	4543	2901
		0.7431	0.9467	0.6552	0.9275	0.1674	0.6526	0.9720	0.7937	0.6830	0.9091
	UvA_ILPS-baseline2	10274	9719	3392	8024	2908	7330	22204	5510	6859	2914
		0.0638	0.0526	0.1921	0.0552	0.8194	0.0869	0.0184	0.1301	0.1221	0.0059
	WaterlooClarke-UWPAH1	307	315	315	315	315	315	315	315	308	308
		0.0553	0.0526	0.1921	0.0552	0.7885	0.0869	0.0180	0.1309	0.1217	0.0059
	WaterlooClarke-UWPAH2	314	315	315	315	315	307	315	308	308	303
		0.9813	0.9988	0.9581	0.9974	0.9163	0.8856	0.9977	0.9878	0.9764	0.0059
	WaterlooCormack-Knee100	15256	11396	5729	10336	1232	15257	27254	6325	8498	130
		0.9830	0.9991	0.9483	0.9976	0.9119	0.8809	0.9978	0.9861	0.9764	0.9802
	WaterlooCormack-Knee1000	16810	11396	5188	10336	1105	15257	27254	5729	8498	3144
		0.9194	0.9979	0.9360	0.9963	0.9692	0.7590	0.9973	0.9857	0.9389	0.9802
	WaterlooCormack-stop2399	7702	10336	4251	9373	2841	5729	24749	5729	5188	3144
		0.7191	0.9639	0.7285	0.7696	0.2819	0.6856	0.9669	0.8863	0.6097	0.0000
	Webis-baseline	5213	7578	1837	4967	231	4044	17786	2938	2554	110
		0.0007	0.9627	0.7174	0.7686	0.3172	0.6875	0.9665	0.8842	0.5912	0.0020
	Webis-keyphrase	105	7351	1817	4858	237	4039	17767	2925	2438	110

Table 7: Set-Based Results (Recall Over Effort) For the athome1 Test Collection. The top number in each cell indicates the recall achieved when the submission indicated by “calling its shot” that a “reasonable” result had been achieved, or when the submission terminated its run. The second number indicates the number of documents submitted to the automated assessment server at this point. (*) indicates the number of additional documents reviewed by a manual review team, and not necessarily submitted to the automated assessment server.

	Topic (R) – Athome2 Collection									
eDiscoveryTeam	2052 (265)	2108 (661)	2129 (589)	2130 (2299)	2134 (252)	2158 (1256)	2225 (182)	2322 (9517)	2333 (4805)	2461 (179)
	0.9698	0.8850	0.9847	0.8534	0.9563	0.9881	0.8956	0.8512	0.8745	0.9777
	871 +2325*	1505 +2101*	3854 +94*	3802 +285*	3689 +19*	1339 +1335*	201 +205*	12799 +195*	6670 +228*	1918 +32*
NINJA	0.8566	0.6702	0.7793	0.4728	0.6191	0.6393	0.9176	0.1567	0.5290	0.6201
	516	990	1276	1867	534	850	665	2497	3310	728
UvA_ILPS-baseline	0.9962	0.9803	0.9745	0.7747	0.8651	0.9881	0.9561	0.9032	0.9463	0.9888
	4651	5038	4657	6363	5141	4651	4652	13969	8224	4653
UvA_ILPS-baseline2	0.9660	0.9183	0.8693	0.7234	0.6865	0.6226	0.8681	0.8880	0.7832	0.8883
	4672	4746	4736	8472	4722	5704	4715	25613	11549	4692
WaterlooClarke-UWPAH1	0.8491	0.3858	0.4414	0.0761	0.7183	0.2508	0.8407	0.0302	0.0651	0.6872
	307	315	324	308	315	315	308	315	313	320
WaterlooCormack-Knee100	0.9811	0.9758	0.9847	0.9622	0.9405	0.0892	0.9231	0.9841	0.9881	0.3799
	989	3144	3478	20402	2316	130	630	30010	15257	130
WaterlooCormack-Knee1000	0.9849	0.9788	0.9796	0.9604	0.9405	0.9881	0.9560	0.9862	0.9875	0.9385
	1104	3144	3144	20402	2316	2841	1105	33042	15257	1372
WaterlooCormack-stop2399	0.9925	0.9818	0.9830	0.7086	0.9524	0.9881	0.9835	0.8800	0.9534	0.9888
	2840	3478	3144	4697	2841	4251	2841	12562	8498	2841
Webis-baseline	0.1094	0.9047	0.0459	0.0000	0.0159	0.6600	0.1539	0.7376	0.7378	0.5475
	173	3002	110	80	111	1457	123	11313	5611	271
Webis-keyphrase	0.0868	0.8684	0.0492	0.0000	0.0159	0.6147	0.1044	0.7354	0.7590	0.4078
	172	3096	110	111	111	1229	111	11152	5915	278

Table 8: Set-Based Results (Recall Over Effort) For the athome2 Test Collection. The top number in each cell indicates recall achieved when the submission indicated by “calling its shot” result had been achieved, or when the submission terminated. The second number indicates the number of documents submitted to the automated assessment server at this point. (*) indicates the number of additional documents reviewed by a manual review team, and not necessarily submitted to the automated assessment server.

	Topic (R) – Athome3 Collection									
eDiscoveryTeam	3089 (255)	3133 (113)	3226 (2094)	3290 (26)	3357 (629)	3378 (66)	3423 (76)	3431 (1111)	3481 (2036)	3484 (23)
	0.9255	0.7699	0.9842	0.8846	0.9173	0.8939	0.4474	0.9973	0.9646	1.0000
	250	97	5347	95	701	106	40	1121	2077	23
	+834*	+49*	+18*	+306*	+920*	+200*	+92*	+272*	+367*	+73*
NINJA	0.7765	0.7699	0.1395	0.8846	0.9603	0.8485	0.4605	0.5995	0.6027	1.0000
	201	143	474	92	714	99	45	673	1392	23
UvA.ILPS-baseline	0.9961	1.0000	0.9895	1.0000	0.9857	0.9546	0.8290	0.9991	0.9509	1.0000
	9024	9024	9031	9024	9100	9024	9936	9024	9026	9024
UvA.ILPS-baseline2	0.9765	0.9735	0.9231	0.6923	0.9142	0.9546	0.4474	0.9937	0.8144	1.0000
	9057	9097	9032	9106	9103	9091	9116	9046	9325	9084
WaterlooClarke-UWPAH1	0.8510	0.9912	0.1423	0.7692	0.4706	0.9848	0.4737	0.2844	0.1552	1.0000
	314	330	315	309	315	332	309	315	315	330
WaterlooCormack-Knee100	0.4353	0.9912	0.9790	0.6923	0.9523	0.8636	0.4605	0.9847	0.9514	1.0000
	129	343	4251	111	1526	130	151	1232	5188	111
WaterlooCormack-Knee1000	0.9961	0.9912	0.9733	1.0000	0.9634	1.0000	0.5263	0.9892	0.9504	1.0000
	1104	1105	3846	1105	1883	1105	1105	1232	4251	1105
WaterlooCormack-stop2399	0.9961	1.0000	0.9852	1.0000	0.9841	1.0000	0.6184	0.9991	0.9504	1.0000
	2840	2566	5188	2566	3144	2566	2566	3846	5188	2566
Webis-baseline	0.6980	0.7788	0.1777	0.5385	0.8935	0.7879	0.3816	0.9829	0.6842	1.0000
	2156	261	942	119	5338	314	130	9971	2215	110
Webis-keyphrase	0.8275	0.0620	0.1380	0.2692	0.8347	0.5000	0.1974	0.9748	0.6680	0.8696
	1954	106	1066	111	5212	296	111	10175	6235	110

Table 9: Set-based Results (recall over effort) for the athome3 test collection. The top number in each cell indicates the recall achieved when the submission indicated by “calling its shot” result had been achieved, or when the submission terminated its run. The second number indicates the number of documents submitted to the automated assessment server at this point. (*) indicates the number of additional documents reviewed by a manual review team, and not necessarily submitted to the automated assessment server.

Topic (R) – Kaime Collection				
	Open	Restricted	Record	VA Tech
	(131698)	(14341)	(166118)	(20083)
UvA.IILPS-baseline	0.9237	0.8458	0.9599	0.9558
	183797	19616	228712	27456
WaterlooClarke-UWPAH1	0.0020	0.0191	0.0013	0.0126
	304	304	304	304
WaterlooCormack-stop2399	0.6605	0.7440	0.8801	0.9507
	104421	15259	185276	27256
Webis-baseline	0.7291	0.7171	0.8858	0.7934
	134983	15523	205524	20854
Webis-keyphrase	0.0001	0.7214	0.9007	0.7916
	73	15806	210913	20785

Table 10: Set-Based Results (Recall Over Effort) For the Kaime Test Collection. The top number in each cell indicates the recall achieved when the submission indicated by “calling its shot” that a “reasonable” result had been achieved, or when the submission terminated its run. The second number indicates the number of documents submitted to the automated assessment server at this point.

	Topic (R) – MIMIC II Collection [Part I]									
	C01	C02	C03	C04	C05	C06	C07	C08	C09	C10
UvA.IILPS-baseline	(5881)	(3867)	(15101)	(7826)	(6123)	(5081)	(19182)	(11256)	(8706)	(8741)
	0.6910	0.7210	0.0031	0.6360	0.5692	0.4718	0.9585	0.0022	0.0030	0.6608
	5960	3528	319	7520	5061	4146	20584	328	317	7520
UvA.IILPS-baseline2	0.5581	0.6289	0.9364	0.6254	0.4023	0.3340	0.9766	0.7831	0.6230	0.6984
	5702	4332	20458	8904	4569	3815	22203	12381	8459	9989
WaterlooClarke-UWPAH1	0.0452	0.0706	0.0195	0.0235	0.0483	0.0382	0.0157	0.0263	0.0326	0.0257
	301	302	302	302	302	302	302	302	302	302
WaterlooCormack-stop2399	0.8109	0.8836	0.9548	0.7313	0.7013	0.6215	0.9949	0.9087	0.8256	0.8396
	8497	6981	20402	9373	7703	6325	27254	15257	11396	11396
Webis-baseline	0.6817	0.7277	0.9533	0.6606	0.5078	0.4399	0.9722	0.8047	0.6149	0.7603
	6630	3709	21110	9082	5098	4360	22035	12790	7580	10727
Webis-keyphrase	0.6667	0.7199	0.9464	0.0017	0.5166	0.0010	0.9691	0.8189	0.6347	0.0008
	6387	3609	20725	47	5283	48	21848	13195	7973	48

Table 11: Set-Based Results (Recall Over Effort) for the MIMIC II Test Collection, Part I. The top number in each cell indicates the recall achieved when the submission indicated by “calling its shot” that a “reasonable” result had been achieved, or when the submission terminated its run. The second number indicates the number of documents submitted to the automated assessment server at this point.

	Topic (R) – MIMIC II Collection [Part II]								
	C11	C12	C13	C14	C15	C16	C17	C18	C19
	(180)	(2579)	(3465)	(2143)	(5143)	(8047)	(11117)	(16827)	(6828)
UvA.ILPS-baseline	0.9500	0.0016	0.0026	0.5684	0.9117	0.4977	0.6980	0.4561	0.5053
	359	315	315	2019	5481	6710	10436	7966	4319
UvA.ILPS-baseline2	0.6778	0.2404	0.0205	0.3985	0.8707	0.3759	0.7603	0.7688	0.4782
	683	1650	388	1569	5249	5699	13813	18688	4979
WaterlooClarke-UWPAH1	0.9444	0.0733	0.0390	0.0854	0.0428	0.0229	0.0247	0.0180	0.0395
	302	302	302	302	302	302	302	302	302
WaterlooCormack-stop2399	0.9889	0.6758	0.5957	0.7485	0.9994	0.5469	0.8239	0.8585	0.6995
	2841	4697	5188	4697	9373	7703	13845	20402	8498
Webis-baseline	0.9556	0.4362	0.3417	0.6122	0.8991	0.5144	0.7857	0.8271	0.5082
	797	2594	2761	2861	5264	7753	13487	21000	4449
Webis-keyphrase	0.8333	0.4017	0.0023	0.5469	0.9078	0.5211	0.8093	0.8126	0.5069
	301	2390	47	2484	5380	7880	14073	20449	4416

Table 12: Set-Based Results (Recall Over Effort) for the MIMIC II Test Collection, Part II. The top number in each cell indicates the recall achieved when the submission indicated by “calling its shot” that a “reasonable” result had been achieved, or when the submission terminated its run. The second number indicates the number of documents submitted to the automated assessment server at this point.

6 Conclusions

The 2015 Total Recall Track successfully deployed a new evaluation architecture, using five new datasets. Three of the datasets will be publicly available, while two datasets must remain private. The results appear to be consistent across all datasets: A fully automated Baseline Model Implementation achieved high recall for all topics, with effort proportionate to the number of relevant documents. Several manual and automatic participant efforts achieved higher recall with less effort than the baseline on some topics, but none consistently improved on the baseline. Some teams appeared to be able to predict when a “reasonable” result had been achieved; however, further work is needed to derive appropriate measures to evaluate what is “reasonable.”

Acknowledgement

We are grateful to the Library of Virginia for affording us access to the Kaine email collection for the Sandbox evaluation. We give particular thanks to Susan Gray Paige, Roger Christman, Rebecca Morgan, and Kathy Jordan for their invaluable assistance with and unwavering support for this work.

References

- [1] Wei Lu Chuan Wu and Ruixue Wang. WHU at TREC Total Recall Track 2015. In *Proc. TREC-2015*, 2015.
- [2] Gordon V. Cormack and Maura R. Grossman. The Grossman-Cormack Glossary of Technology-Assisted Review. *Fed. Cts. L. Rev.*, 7(1), 2013.
- [3] Gordon V Cormack and Maura R Grossman. Multi-faceted recall of continuous active learning for technology-assisted review. In *Proceedings of the 38th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 763–766. ACM, 2015.
- [4] Gordon V. Cormack and Maura R. Grossman. Waterloo (Cormack) Participation in the TREC 2015 Total Recall Track. In *Proc. TREC-2015*, 2015.
- [5] Gordon V Cormack and Mona Mojdeh. Machine Learning for Information Retrieval: TREC 2009 web, relevance feedback and legal tracks. In *Proc. TREC-2009*, 2009.
- [6] Evangelos Kanoulas David van Dijk, Zhaochun Ren and Maarten de Rijke. The University of Amsterdam (ILPS) at TREC 2015 Total Recall Track. In *Proc. TREC-2015*, 2015.
- [7] Vivek Dhand. Efficient semantic indexing via neural networks with dynamic supervised feedback. In *Proc. TREC-2015*, 2015.
- [8] Maura R Grossman and Gordon V Cormack. Comments on “The implications of Rule 26(g) on the use of technology-assisted review”. *Fed. Cts. L. Rev.*, 8(1), 2014.
- [9] Yipeng Wang Charles L.A. Clarke Haotian Zhang, Wu Lin and Mark D. Smucker. WaterlooClarke: TREC 2015 Total Recall Track. In *Proc. TREC-2015*, 2015.
- [10] Bruce Hedin, Stephen Tomlinson, Jason R. Baron, and Douglas W. Oard. Overview of the TREC 2009 Legal Track. In *Proc. TREC-2009*, 2009.
- [11] Julian PT Higgins, Sally Green, et al. *Cochrane handbook for systematic reviews of interventions*, volume 5. Wiley Online Library, 2008.
- [12] Bayu Hardi Jeremy Pickens, Tom Gricks and Mark Noel. A Constrained Approach to Manual Total Recall. In *Proc. TREC-2015*, 2015.
- [13] Ralph Losey, Jim Sullivan, and Tony Reichenberger. e-Discovery Team at TREC 2015 Total Recall Track. In *Proc. TREC-2015*, 2015.
- [14] Mihai Lupu. TUW at the first Total Recall Track. In *Proc. TREC-2015*, 2015.
- [15] Magdalena Keil Olaoluwa Anifowose Amir Othman Matthias Hagen, Steve Göring and Benno Stein. Webis at TREC 2015: Tasks and Total Recall Tracks. In *Proc. TREC-2015*, 2015.

- [16] Adam Roegiest and Gordon V Cormack. Total Recall Track Tools Architecture Overview. In *Proc. TREC-2015*, 2015.
- [17] Hui Yang, John Frank, and Ian Soboroff. Trec 2015 dynamic domain track overview. In *Proc. TREC-2015*, 2015.