

# Cluster-Based Adaptive and Batch Filtering

David Eichmann, Miguel Ruiz, Padmini Srinivasan

School of Library and Information Science

University of Iowa

## 1 – Introduction

Information filtering is increasingly critical to knowledge workers drowning in a growing flood of byte streams [6, 8, 9]. Our interest in the filtering track for TREC-7 grew out of work originally designed for information retrieval on the Web, using both ‘traditional’ search engine [5] and agent-based techniques [6, 7]. In particular, the work by Cutter, et. al. in clustering [3, 4] has great appeal in the potential for synergistic interaction between user and retrieval system.

Our efforts for TREC-7 included two distinct filtering architectures, as well as a more traditional approach to the adhoc track (which used SMART 11.3). The filtering work was done using TRECCer – our Java-based clustering environment – alone for adaptive filtering and a combination of TRECCer and SMART for batch filtering.

## 2 – Adhoc Track

### 2.1 – Adhoc Methodology

Our overall approach is to apply Rocchio-based retrieval feedback [11] for query expansion. The second run submitted (Iowacuhk2) is simply such a run with the top 10 documents from an initial retrieval run assumed relevant and the best 350 terms extracted from these documents. Documents and queries were weighted using the Lnu.ltu scheme [12] which had yielded good results in previous TREC runs [e.g., 2]. For our primary run (Iowacuhk1), we focussed on improving the initial retrieved set that is assumed relevant during retrieval feedback. The following steps describe our approach.

1. Retrieve 10 documents using the initial query (Lnu.ltu) weights. Call this set A.
2. Identify the top 3 documents for each query.
3. Treat these top 3 documents as pseudo-queries, index them against the same database and retrieve 100 documents for each pseudo query (Lnu.ltu weights with pivot average document size of 126 and a slope of 0.2).
4. Merge the 100 documents from the three pseudo-queries and eliminate duplicates. Call this set B.
5. Find the intersection of sets A and B for each topic. Use this set for retrieval feedback to expand the query. Call this set C.
6. Expand the original query by 350 terms using  $\alpha = 8$ ,  $\beta = 8$  and  $\gamma = 8$  with set C and using Rocchio's algorithm.
7. Retrieve the final set of 1000 documents using the expanded query (Lnu.ltu weights using the above parameters).

### 2.2 – Results and Analysis

Figures 1 through 3 compare the performance of the two Iowa runs against the minimum, maximum and median scores for the adhoc track. Table 1 below summarizes this performance. It shows the number of topics in which the corresponding Iowa run performs at or better than the median value.

**Table 1: Adhoc Performance**

	$\geq$ Median, Top 100 Retrieved	$\geq$ Median, Top 1000 Retrieved	$\geq$ Median, Avg. Precision	Avg. Precision (non-interpolated)	Exact Precision
Iowacuhk1	37	32	32	0.2221	0.2680
Iowacuhk2	39	34	35	0.2260	0.2754

Thus in 64 to 78% of the topics, the Iowa runs are at or above the median performance. It is also seen that our second run, i.e., the straight Lnu.ltu and Rocchio-based retrieval feedback approach is slightly better for each measure than our primary run in which we tried to refine the set of documents used for retrieval feedback. Although somewhat

## Cluster-Based Adaptive and Batch Filtering

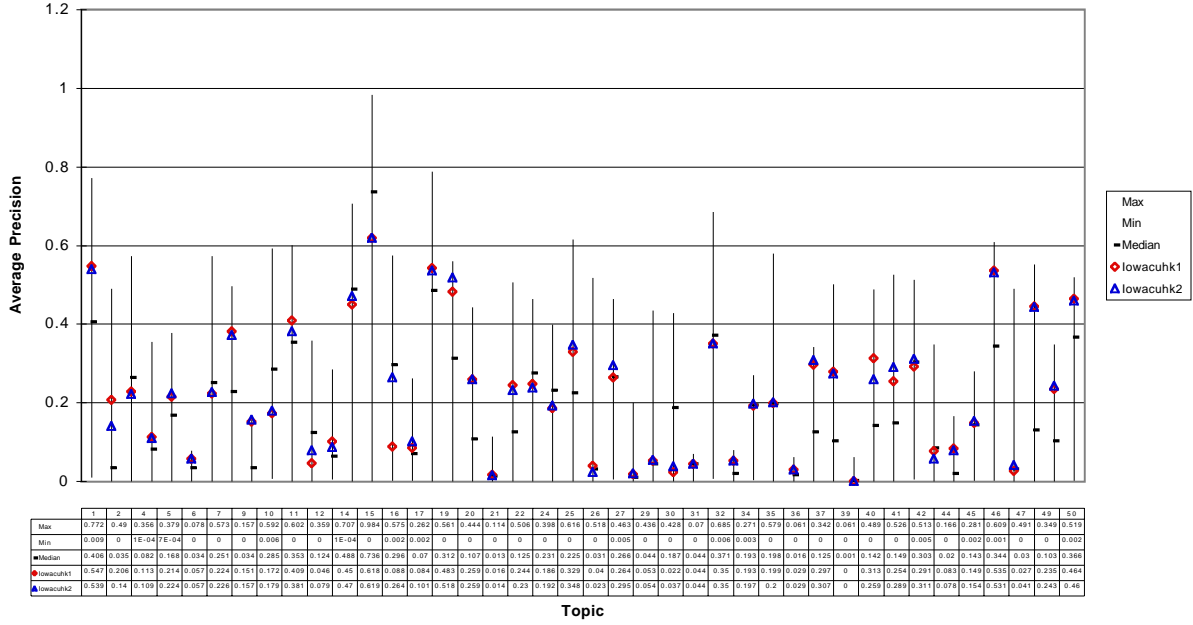


Figure 1: Adhoc Retrieval, Average Precision

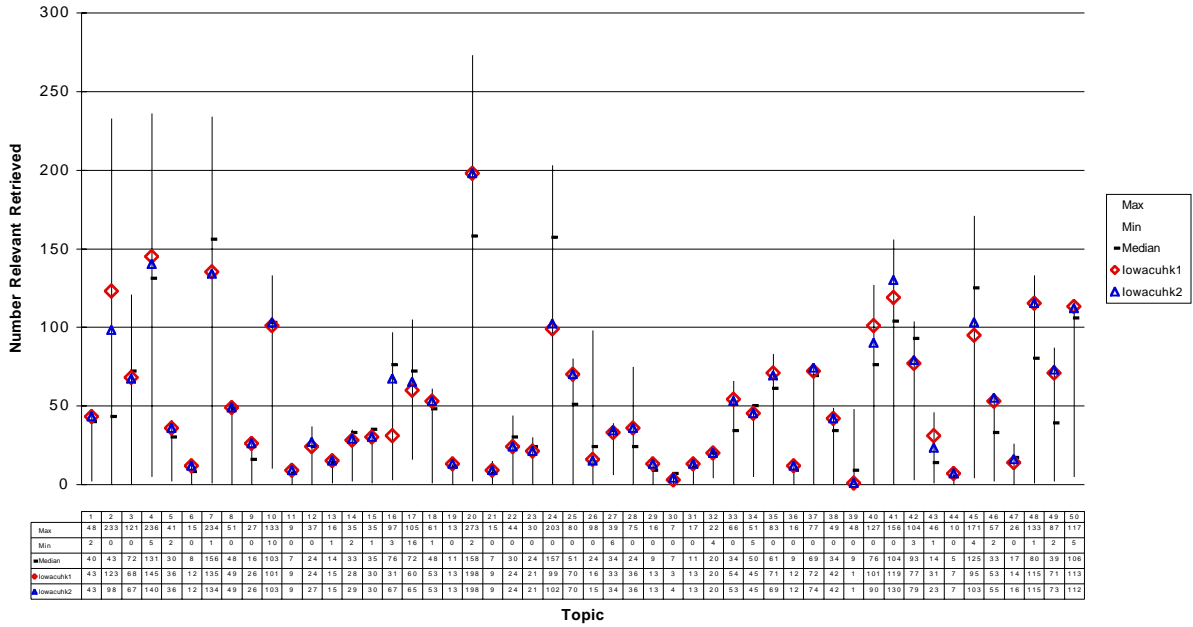


Figure 2: Adhoc Retrieval, # Relevant Retrieved in Top 1000 Documents

disappointing, this performance sets an internal baseline against which we hope to show improvements in our future TREC efforts.

### 3 – Adaptive Filtering Track

Our existing approach to Web search/filtering involves a dynamic clustering technique where the threshold for formation of new clusters and the threshold for visibility of ‘sufficiently important’ clusters can be specified by the

## Cluster-Based Adaptive and Batch Filtering

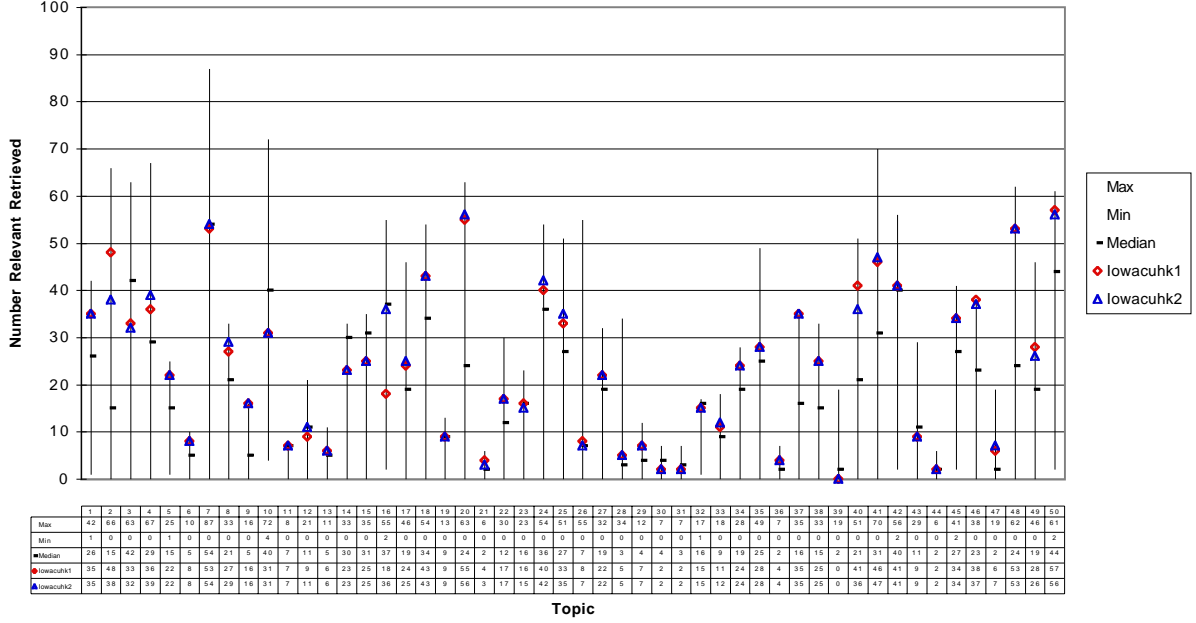


Figure 3: Adhoc Retrieval, # Relevant Retrieved in Top 100 Documents

user when the topic is presented to the system. As documents are retrieved and clusters form, the user can select interesting clusters for further exploration through the links contained in member documents [7]. The TREC requirements for multi-query support and simulation of user judgment responses led us to modify the single set-of-clusters model, creating a two-level scheme.

Similarity between documents and clusters is measured using a straight-forward vector cosine measure:

$$sim(d, c) = \frac{\sum_{i=1}^{N_d} \sum_{j=1}^{N_c} (TF(W_i) \cdot TF(W_j))}{\sqrt{\sum_{i=1}^{N_d} TF(W_i) \cdot \sum_{j=1}^{N_c} TF(W_j)}}$$

Term frequencies are built up incrementally as a given run progresses and cluster term weights are adjusted every ten input files. This approach is therefore somewhat inaccurate in the initial phases of a run, but quickly reaches a point of reasonable stability with respect to term frequencies and has the added benefit of requiring no fore-knowledge of the vocabulary. All vocabulary is stemmed using Porter's algorithm [10]. We prune document term vectors to the 100 most weighty terms and cluster vectors to the 200 most weighty terms. This proves to have no significant effect on the accuracy of our results, but a significant effect on both memory requirements and execution time, the latter due to a corresponding reduction in the cost of dot product calculations.

The primary cluster level corresponds to the internal representation of a topic definition. The original threshold specifications were retained here to allow specification of the first-order recall of the system. We experimented with a variety of means of generating a primary similarity measure, but settled on one based upon the text of the topic's concept definitions for the submitted runs.

The secondary level is where the adaptive portion of the system functions and where we found the most benefit in parameter tuning. Each primary cluster (and hence, each topic) has a private set of zero or more secondary clusters. When a document clears the threshold for a primary cluster, it either joins an existing secondary cluster or forms a new one, based upon a membership threshold. The shift from a single membership threshold to a primary/secondary pair

## Cluster-Based Adaptive and Batch Filtering

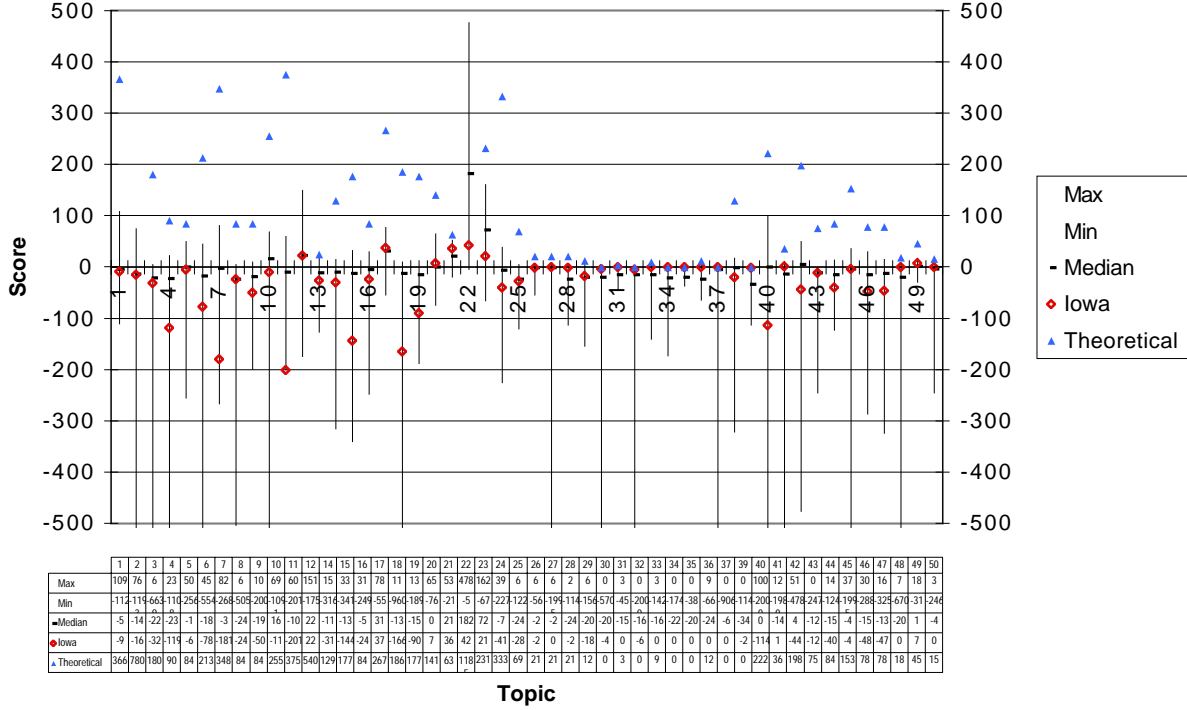


Figure 4: Adaptive Filtering, AP'88 F1 Results

allowed us to achieve a tunable level of recall (by using a lower primary threshold, as mentioned above) while teasing out distinctions between candidate document clusters through use of a higher secondary threshold.

Introduction of a visibility threshold for secondary cluster similarity to the primary then gives us a means for adaptation. When a secondary cluster's similarity first exceeds the visibility threshold, its member documents are declared to the user and relevance judgments are obtained. The secondary is then colored appropriately. Secondary clusters containing relevant (and unjudged, if any) documents are colored green and have all subsequent members declared as relevant. Secondaries containing non-relevant (and potentially, unjudged) documents are colored red and declare no further members. Adaptation then occurs over time as secondary clusters exceed the visibility threshold and are colored, with red secondary clusters mitigating the lack of precision provided by the recall-centric primary threshold.

Secondary clusters exceeding the visibility threshold potentially contain a mix of different document types (relevant, non-relevant and unjudged). We currently address this in the following, conservative manner: if a secondary cluster contains

- one or more relevant documents, no non-relevant documents and zero or more unjudged documents, color it green;
- one or more non-relevant documents, no relevant documents and zero or more unjudged documents, color it red;
- one or more relevant documents, one or more non-relevant documents and zero or more unjudged documents, color it gray, but treat it as green;
- fewer than a specific number (currently 10) of unjudged documents and no relevant or non-relevant documents, leave it uncolored until the first relevant or non-relevant document is added, then color it appropriately (note that this optimistic stance has a distinct effect w.r.t. false positives); and finally,
- more than a specific number of unjudged documents and no relevant or non-relevant documents, color it red (we do this pessimistically due to the low density of judged documents in the corpus).

We selected a primary similarity threshold of 0.18, secondary similarity threshold of 0.5 and visibility threshold of 3 in our preliminary experiments with the Wall Street Journal corpus, but used a primary similarity threshold of 0.15, secondary similarity threshold of 0.4 and visibility threshold of 2 based upon an assumption that the WSJ corpus involved a more restricted vocabulary than the AP vocabulary. Figures 4, 5, and 6 show our results for AP88, AP89,

## Cluster-Based Adaptive and Batch Filtering

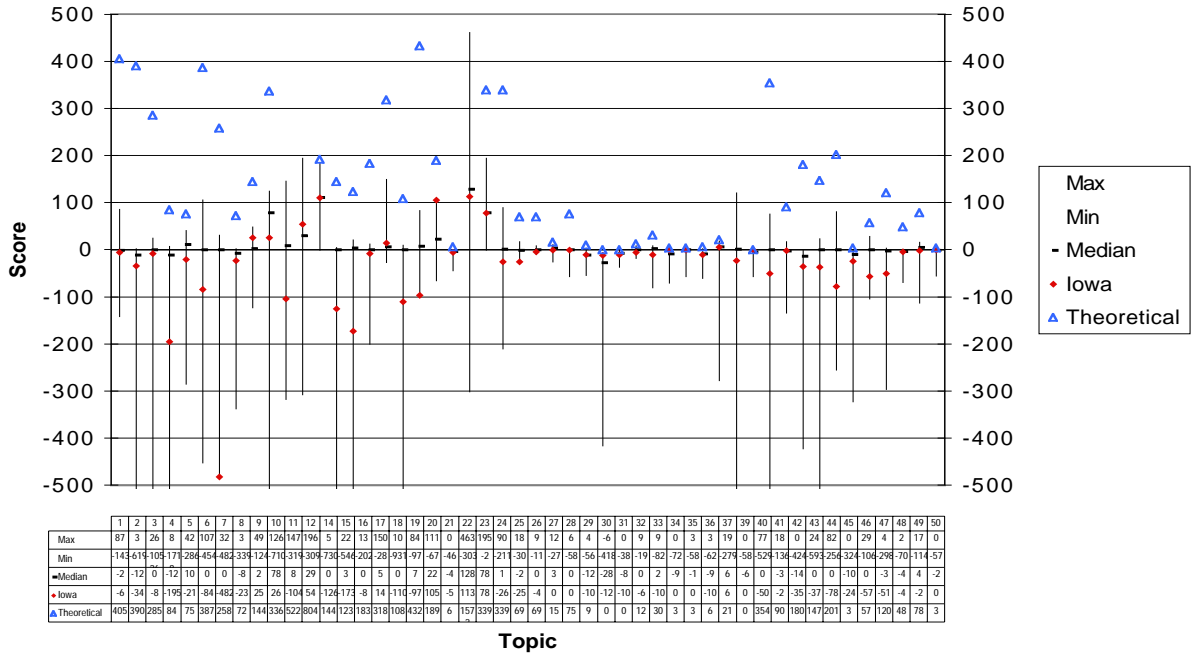


Figure 5: Adaptive Filtering, AP'89 F1 Results

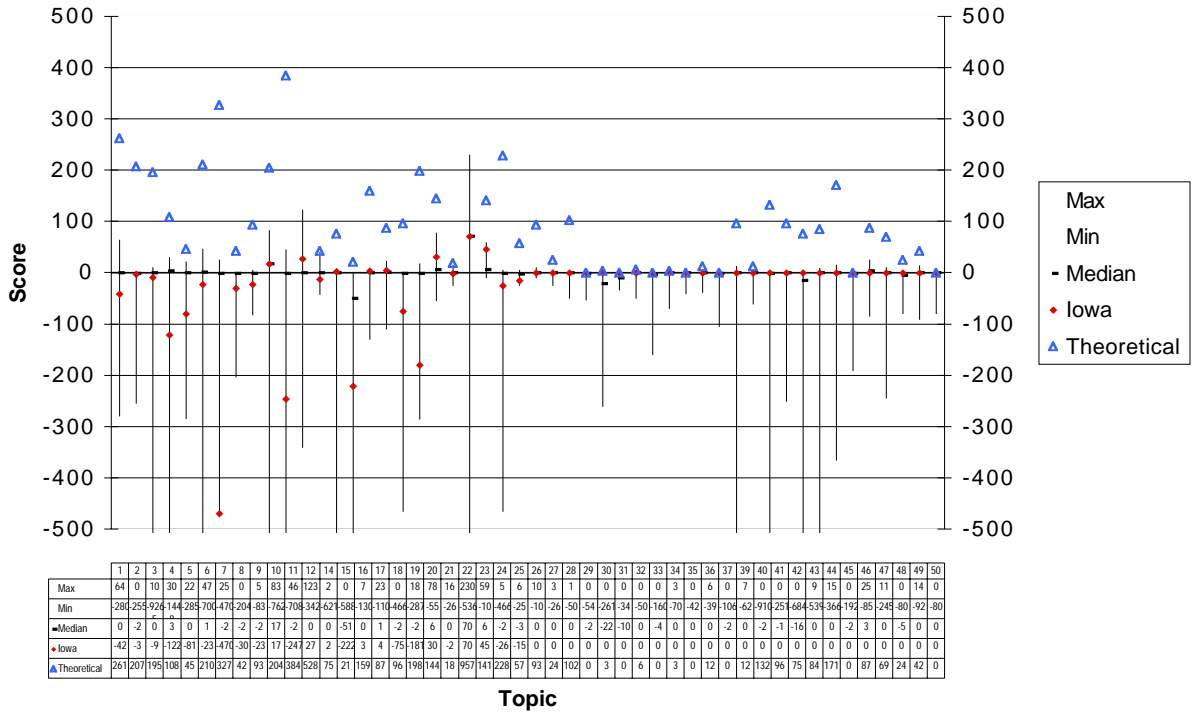


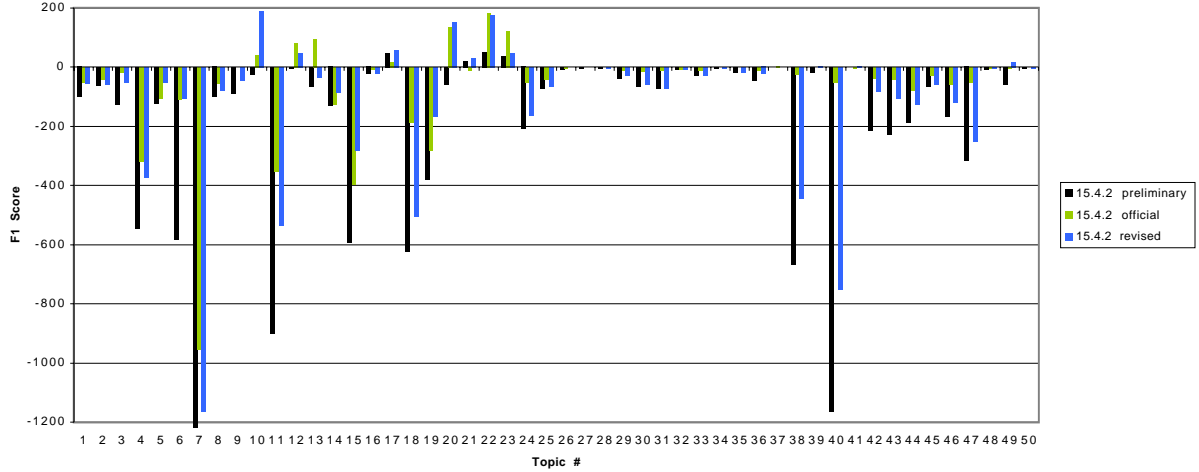
Figure 6: Adaptive Filtering, AP'90 F1 Results

and AP90, respectively for the F1 measure. The vertical bars indicate the min, max and median scores for a given topic, the circles our score, and the diamonds the theoretically possible score. Table 2 shows the number of topics for which we score at or above the median for each of the three years.

## Cluster-Based Adaptive and Batch Filtering

**Table 2: Adaptive Filtering Result Comparison**

Year	# of scores above median
AP88	23
AP89	15
AP90	31



**Figure 7: F1 Scoring Shifts Due to QRel Coverage**

With a few notable exceptions (e.g., Topic 7), our approach scores very near the best performance recorded in the official runs. This includes the topics with extremely low density of relevant documents. Our performance in suppressing the negatively judged documents in topics 26-50 appears to train well in ‘88-’89, but at the cost of missing relevant documents in ‘90 in topics 40-50. This could well be an artifact of the trade-off between the number of secondary clusters formed and their resulting specificity. A smaller number of more general clusters suffers from lack of focus, but also declares fewer negatively judged documents in coloring a cluster. A larger number of more specific clusters improves cluster specificity, but can cause the declaration of a greater number of negatively judged documents as clusters corresponding to fine discriminations in concepts present in the corpus are colored. These effects will require substantially more experimentation.

The points in time at which feedback occurs has significant effect upon the performance of the clustering algorithm. Figure 7 shows the results for our preliminary run (with original qrels for ‘88 and ‘89 only), the official run (with original qrels for all years) and a revised run (with revised qrels for all years). In virtually all topics the cumulative F1 score is higher for both the official runs and the revised qrel run over the preliminary run with qrels for only ‘88 and ‘89. Note however, that there are numerous cases where the official run outperforms the revised qrel run. We suspect that this is due to shifts in cluster declaration patterns across the changing qrel patterns. As an example of this, consider the pattern of black cluster declaration for Topic 5 as shown in Figures 8 and 9. As the preliminary run exhausts its qrels at snapshot 69 (the end of ‘89), no additional black clusters are declared until the end of ‘90. The revised run, with additional judgments for ‘90, continues to declare a substantial number of negative clusters, but at the same time, scores better than the preliminary run. The growth in the number of negative clusters is due to the relatively high density of negative judgements compared to positive judgements in ‘90.

Figures 10 and 11 show a somewhat different situation for Topic 12. While the density of all judgements varies more substantially in ‘90, the number of both positive and negative clusters declared does not substantially differ from the preliminary run to the revised run. The cumulative F1 score, however, shows a distinct improvement, capitalizing on already declared clusters to suppress non-relevant documents and declare relevant documents.

## Cluster-Based Adaptive and Batch Filtering

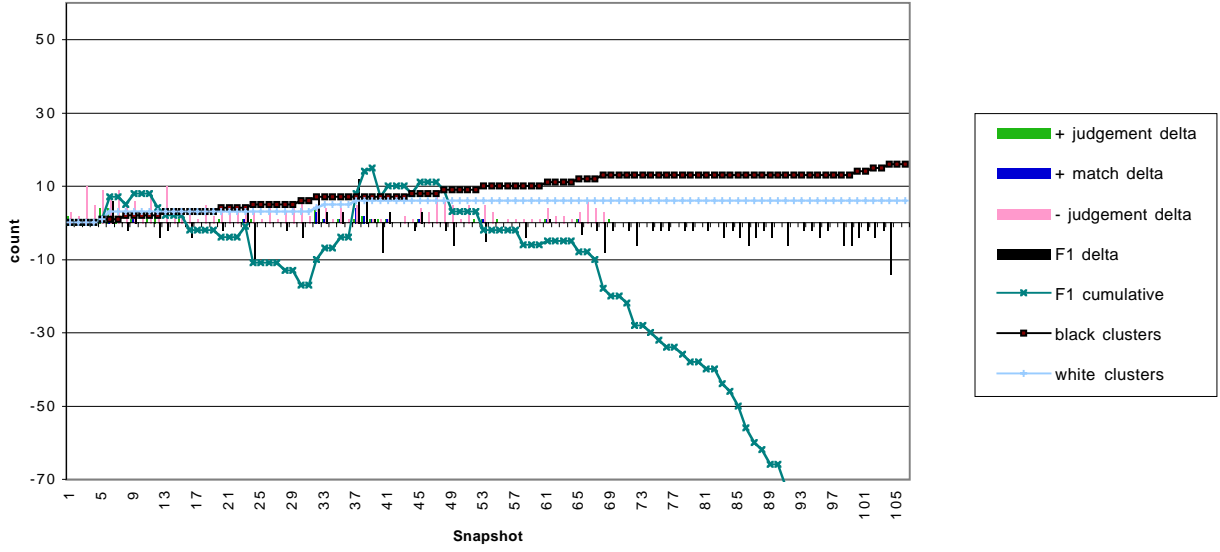


Figure 8: Topic 5, Preliminary QRels

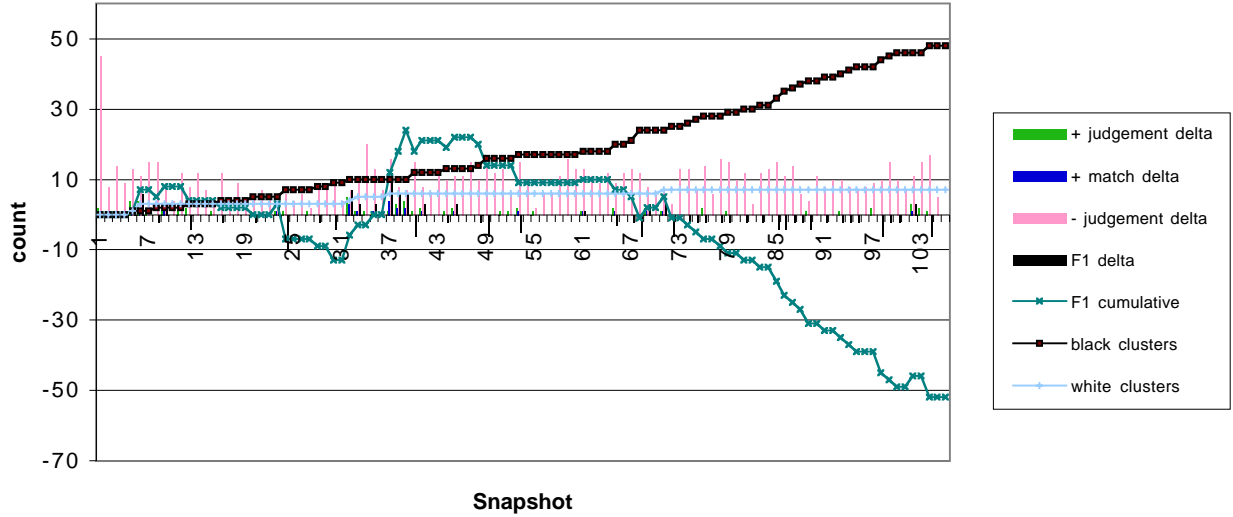


Figure 9: Topic 5, Revised QRels

### 4 – Batch Filtering Track

For this subtrack we decided to use Rocchio's query expansion method to build an initial profile for each topic, and then let the adaptive system learn on the training before processing the test set. The following steps describe our methodology:

1. Index the training set (AP88, topics, and relevance judgments of AP88) using SMART. We used the pivoted normalization weighting scheme (Lnu.ltu), stop words, and no stemming.
2. Expand the topics using relevance feedback on the training dataset. For this, the initial retrieval run extracted the top 100 documents. Rocchio's method was used with parameters  $\alpha=8$ ,  $\beta=8$ , and  $\gamma=8$ . The top 200 terms were used for expansion. These expanded topics were input to the TRECcer program in step 3.
3. Run the TRECcer program. Originally we intended to generate IDF statistics on all the AP88 train set for the TRECcer program. However, due to lack of time we selected a subset of files and used only those files for IDF statistics and initial cluster formation with the TRECcer. For each topic, the first file in the training collection that

## Cluster-Based Adaptive and Batch Filtering

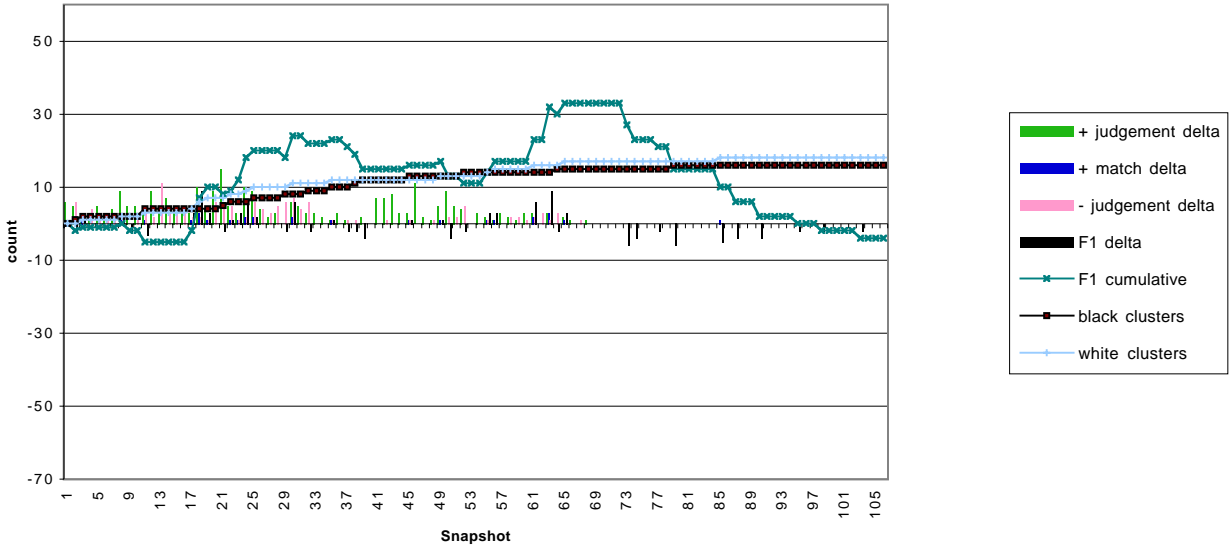


Figure 10: Topic 12, Preliminary QRels

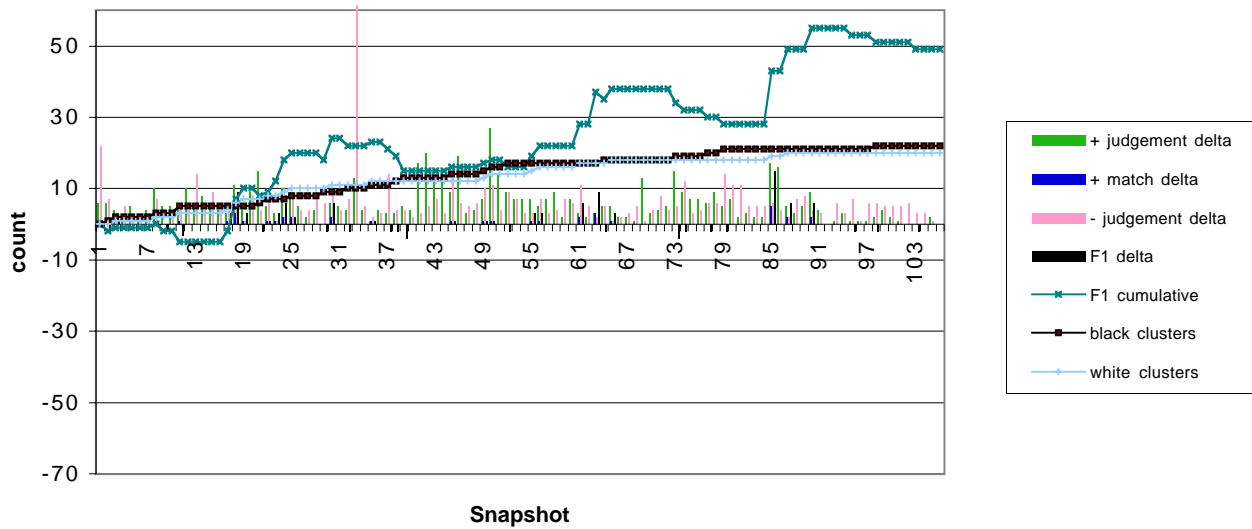


Figure 11: Topic 12, Revised QRels

contains a relevant document was identified. This procedure generated a subset of 26 files.

The parameters used in the batch filtering TRECcr runs were tuned using the training set AP88. We obtained the following settings:

- F1 run (shown in Figure 12): primary cluster membership= 0.25, secondary membership threshold= 0.5, visibility threshold = 0.25.
- F3 run (shown in Figure 13): primary cluster membership= 0.2, secondary membership threshold= 0.5, visibility threshold = 0.25.

We ran the TRECcr program on 20 machines (a mix of HP and SGI workstations), assigning 5 topics to each.

### 4.1 – Results

Unfortunately, the official results included partial results for many of the queries (only queries 36-50 for the F1 run had been completed). We did complete the runs later using the entire AP88 subset for training the system before starting to process the test set (AP89-90).



## Cluster-Based Adaptive and Batch Filtering

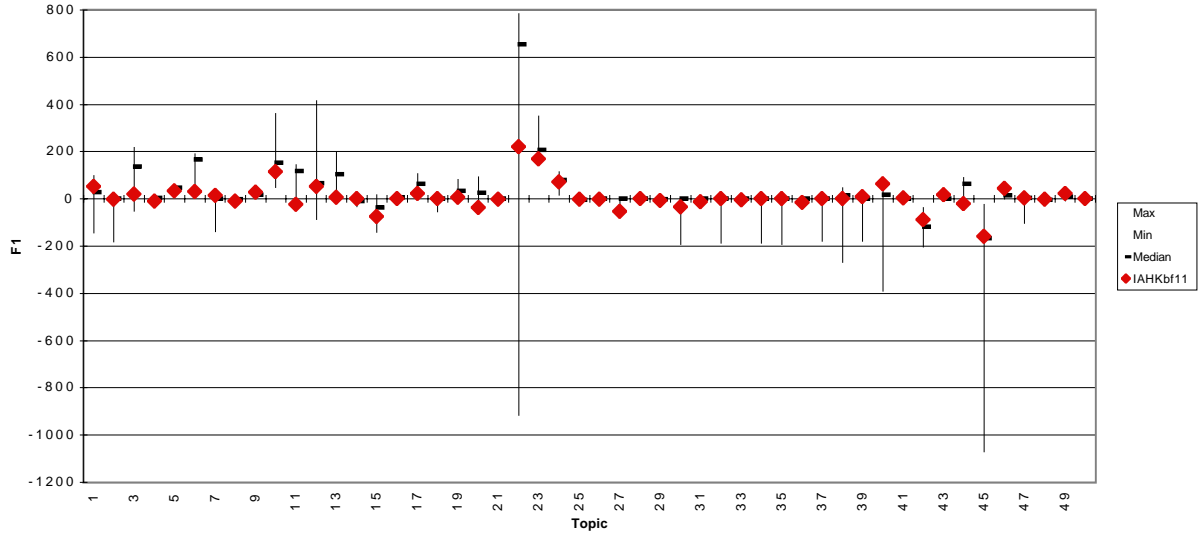


Figure 12: Batch Filtering, F1 Results

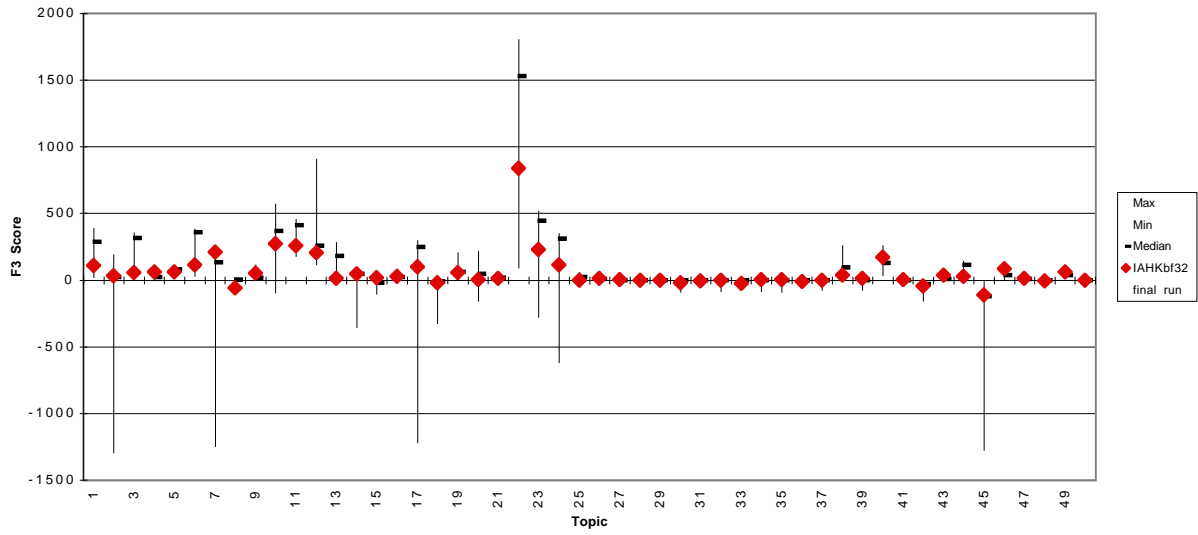


Figure 13: Batch Filtering, F3 Results

Table 3 summarizes the results of our F1 runs. In the official run 17 of the topics were above or equal to the median(12 of them the maximum), while in the unofficial but complete run this improved to 23 (13 of the maximum).

Table 3: Performance of the Batch Filtering Runs

	$\geq$ Median	Maximum in
F1 (official)	17	12
F1 (complete)	22	13
F3 (official)	12	7
F3 (complete)	23	7

## Cluster-Based Adaptive and Batch Filtering

In the official F3 runs 12 of the topics were above the median (7 of them the maximum), whereas in the unofficial run 23 topics were above the median (7 of the maximum).

The results obtained in our F3 runs indicate that the unofficial run significantly outperforms our official run. This improvement is justified in part by the fact that the unofficial runs include more declared documents that eventually will increase the chance of retrieving relevant documents.

In the case of the F1 runs the unofficial run is better but the improvement is smaller. We attribute this to the fact that our official F1 run already had 15 completed topics. There is also a difference in the training examples used. The Official run uses a subset of 26 files, while the unofficial run uses the entire training set. We compared the results obtained in the subset of queries 36-40. In the official run 10 of these topics are above the median (6 at the maximum), while in the unofficial run 12 of the topics are above the median (8 at the maximum). We also observe that there are differences in the number of declarations – 7 topics increase the number of declarations while 2 decrease. This is because a different training set induces a different secondary cluster structure.

## 5 – Conclusions and Future Plans

Our preliminary experience with two-level clustering and a mixed architecture of TRECCer and SMART have been encouraging. We expect that with further tuning of primary/secondary cluster interaction we will achieve significantly better results. Our performance on the Wall Street Journal corpus during our earlier experiments with clustering lead us to believe that similarity thresholds are sensitive to vocabulary diversity, particularly compared to the more diverse vocabulary of the AP corpus. We are quite interested in exploring a blend of lower primary thresholds and higher secondary thresholds. This should improve our recall, but only at the cost of early training of negative examples.

## 6 – References

- [1] Baclace, P. E., “Competitive Agents for Information Filtering,” *Communications of the ACM*, v. 35, n. 12, December 1992, p. 50.
- [2] Buckley, C., M. Mitra, J. Walz, and C. Cardie, “Using Clustering and SuperConcepts within SMART,” *Proc. of Sixth Text REtrieval Conference (TREC-6)*, Gaithersburg, Maryland, November 19-21, 1997. <http://trec.nist.gov/pubs/trec6/papers/index.track.html>
- [3] Cutting, D. R., D. R. Karger and J. O. Pedersen, “Constant interaction-time scatter/gather browsing of very large document collections,” *Proc. of SIGIR’93*, June 1993.
- [4] Cutting, D., D. Karger, J. Pedersen, and J. W. Tukey, “Scatter/Gather: A Cluster-based Approach to Browsing Large Document Collections,” *Proceedings of the 15th Annual International ACM/SIGIR Conference*, Copenhagen, 1992.
- [5] Eichmann, D., “The RBSE Spider – Balancing Effective Search Against Web Load,” *First International Conference on the World Wide Web*, Geneva, Switzerland, May 25-27, 1994, pages 113-120.
- [6] Eichmann, D., “Search and Metasearch on a Diverse Web,” *First W3C Workshop on Distributed Indexing and Search*, Boston, MA, May 28-29, 1996.
- [7] Eichmann, D., “Ontology-Based Information Fusion,” *Workshop on Real-Time Intelligent User Interfaces for Decision Support and Information Visualization, 1998 International Conference on Intelligent User Interfaces*, San Francisco, CA, January 6-9, 1998.
- [8] Lieberman, H., “Letizia: An Agent That Assists Web Browsing,” *Proceedings of the International Joint Conference on Artificial Intelligence*, Montreal, August 1995.
- [9] Lieberman, H., “Autonomous Interface Agents,” *Proceedings of the ACM Conference on Computers and Human Interfaces (CHI’97)*, Atlanta, GA, March 1997.
- [10] Porter, M. F., “An Algorithm for Suffix Stripping,” *Program*, v. 14, no. 3, 1980, p. 130-137.
- [11] Rocchio, J., “Relevance feedback in information retrieval,” in *The SMART Retrieval System: Experiments in Automatic Document Processing*, G. Salton (ed.), Prentice-Hall, p. 313-323.
- [12] Singhal, A., G. Salton, M. Mitra and C. Buckley, “Document Length Normalization,” *Information Processing & Management*, v. 32, no. 5, Sept. 1996, p. 619-633.