# Manual and Automatic Evaluation of Summaries

Chin-Yew Lin and Eduard Hovy

USC Information Sciences Institute

4676 Admiralty Way

Marina del Rey, CA 90292

+1-310-448-8711/8731

{cyl,hovy}@isi.edu

## Abstract

In this paper we discuss manual and automatic evaluations of summaries using data from the Document Understanding Conference 2001 (DUC-2001). We first show the instability of the manual evaluation. Specifically, the low inter-human agreement indicates that more reference summaries are needed. To investigate the feasibility of automated summary evaluation based on the recent BLEU method from machine translation, we use accumulative n-gram overlap scores between system and human summaries. The initial results provide encouraging correlations with human judgments, based on the Spearman rank-order correlation coefficient. However, relative ranking of systems needs to take into account the instability.

## 1    Introduction

Previous efforts in large-scale evaluation of text summarization include TIPSTER SUMMAC (Mani et al. 1998) and the Document Understanding Conference (DUC) sponsored by the National Institute of Standards and Technology (NIST). DUC aims to compile standard training and test collections that can be shared among researchers and to provide common and large scale evaluations in single and multiple document summarization for their participants.

In this paper we discuss manual and automatic evaluations of summaries using data from the Document Understanding Conference 2001 (DUC-2001). Section 2 gives a brief overview of the evaluation procedure used in DUC-2001 and the Summary Evaluation Environment (SEE) interface used to support the DUC-2001 human evaluation protocol. Section 3 discusses evaluation metrics. Section 4 shows the instability of manual evaluations. Section 5 outlines a method of automatic summary evaluation using accumulative n-gram matching score (NAMS) and proposes a view that casts summary evaluation as a decision making process. It shows that the NAMS method is bounded and in most cases not usable, given only a single reference summary to compare with. Section 6 discusses why this is so, illustrating various forms of mismatching between human and system summaries. We conclude with lessons learned and future directions.

## 2    Document Understanding Conference (DUC)

DUC2001 included three tasks:

- Fully automatic single-document summarization: given a document, participants were required to create a generic 100-word summary. The training set comprised 30 sets of approximately 10 documents each, together with their 100-word human written summaries. The test set comprised 30 unseen documents.

- Fully automatic multi-document summarization: given a set of documents about a single subject, participants were required to create 4 generic summaries of the entire set, containing 50, 100, 200, and 400 words respectively. The document sets were of four types: a single natural disaster event; a single event; multiple instances of a type of event; and information about an individual. The training set comprised 30 sets of approximately 10 documents, each provided with their 50, 100, 200, and 400-word human written summaries. The test set comprised 30 unseen sets.
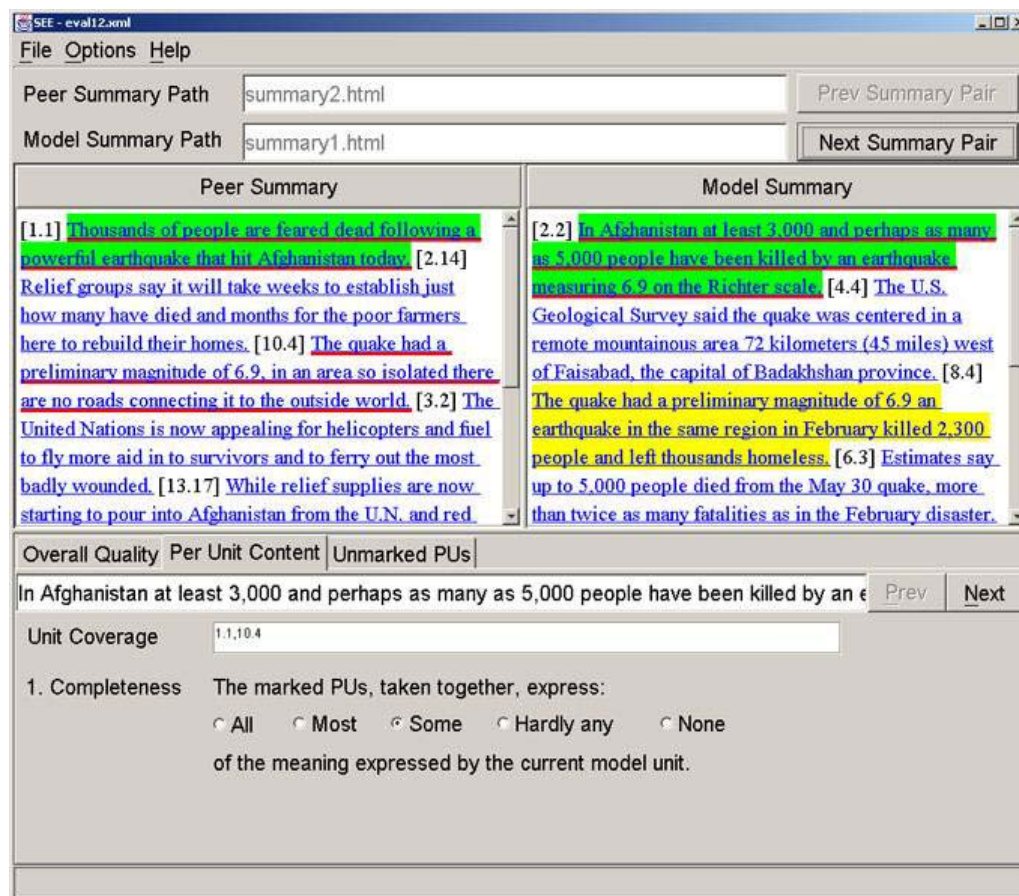
**Figure 1.** SEE in an evaluation session.

- Exploratory summarization: participants were encouraged to investigate alternative approaches to evaluating summarization and report their results.

A total of 11 systems participated in the single-document summarization task and 12 systems participated in the multi-document task.

The training data were distributed in early March of 2001 and the test data were distributed in mid-June of 2001. Results were submitted to NIST for evaluation by July 1st 2001.

## 2.1 Evaluation Materials

For each document or document set, one human summary was created as the 'ideal' model summary at each specified length. Two other human summaries were also created at each length. In addition, baseline summaries were created automatically for each length as reference points. For the multi-document summarization task, one baseline, *lead baseline*,

took the first 50, 100, 200, and 400 words in the last document in the collection. A second baseline, *coverage baseline*, took the first sentence in the first document, the first sentence in the second document and so on until it had a summary of 50, 100, 200, or 400 words. Only one baseline (baseline1) was created for the single document summarization task.

## 2.2 Summary Evaluation Environment

NIST assessors who created the 'ideal' written summaries did pairwise comparisons of their summaries to the system-generated summaries, other assessors' summaries, and baseline summaries. They used the Summary Evaluation Environment (SEE) 2.0 developed by one of the authors (Lin 2001) to support the process. Using SEE, the assessors compared the system's text (the *peer* text) to the ideal (the *model* text). As shown in Figure 1, each text was decomposed into a list of units and displayed in separate windows. In DUC-2001 the sentence was used as the smallest unit of evaluation.

SEE 2.0 provides interfaces for assessors to judge both the content and the quality of summaries. To measure content, assessors step through each model unit, mark all system units sharing content with the current model unit (shown in green highlight in the model summary window), and specify that the marked system units express *all, most, some* or *hardly any* of the content of the current model unit. To measure quality, assessors rate grammaticality[1], cohesion[2], and coherence[3] at five different levels: *all, most, some, hardly any*, or *none*.

For example, as shown in Figure 1, an assessor marked system units 1.1 and 10.4 (shown in red underlines) as sharing *some* content with the current model unit 2.2 (highlighted green).

## 3    Evaluation Metrics

One goal of DUC-2001 was to debug the evaluation procedures and identify stable metrics that could serve as common reference points. NIST did not define any official performance metric in DUC-2001. It released the raw evaluation results to DUC-2001 participants and encouraged them to propose metrics that would help progress the field.

### 3.1    Recall, Coverage, Retention and Weighted Retention

Recall at different compression ratios has been used in summarization research to measure how well an automatic system retains important content of original documents (Mani and Maybury 1999). Assume we have a system summary $S_s$ and a model summary $S_m$. The number of sentences occurring in $S_s$ is $N_s$, the number of sentences in $S_m$ is $N_m$, and the number in both $S_s$ and $S_m$ is $N_a$. Recall is defined as $N_a/N_m$. The Compression Ratio is defined as the length of a summary (by words or sentences) divided by the length of its original document.

Applying this direct all-or-nothing recall in DUC-2001 without modification is not appropriate because:

---

[1] Does the summary observe English grammatical rules independent of its content?

[2] Do sentences in the summary fit in with their surrounding sentences?

[3] Is the content of the summary expressed and organized in an effective way?

1. Multiple system units contribute to multiple model units.

2. Exact overlap between $S_s$ and $S_m$ rarely occurs.

3. Overlap judgment is not binary.

For example in Figure 1, an assessor judged system units 1.1 and 10.4 sharing *some* content with model unit 2.2. Unit 1.1 says "*Thousands of people are feared dead*" and unit 2.2 says "*3,000 and perhaps ... 5,000 people have been killed*". Are "thousands" equivalent to "3,000 to 5,000" or not? Unit 10.4 indicates it was an "*earthquake of magnitude 6.9*" and unit 2.2 says it was "*an earthquake measuring 6.9 on the Richter scale*". Both of them report a "6.9" earthquake. But the second part of system unit 10.4, "*in an area so isolated…*", seems to share some content with model unit 4.4 "*the quake was centered in a remote mountainous area*". Are these two equivalent? This example highlights the difficulty of judging the content coverage of system summaries against model summaries and the inadequacy of using simple recall as defined.

For this reason, NIST assessors not only marked the segments shared between system units (SU) and model units (MU), they also indicated the degree of match, i.e., *all*, *most*, *some*, *hardly any*, or *none*. This enables us to compute *weighted recall*.

Different versions of weighted recall were proposed by DUC-2001 participants. (McKeown et al. 2001) treated the completeness of coverage as a threshold: 4 for *all*, 3 for *most* and above, 2 for *some* and above, and 1 for *hardly any* and above. They then proceeded to compare system performances at different threshold levels. They defined recall at threshold $t$, $Recall_t$, as follows:

$$\frac{\text{Number of MUs marked at or above } t}{\text{Total number of MUs in the model summary}}$$

Instead of thresholds, we use here as coverage score the ratio of completeness of coverage $C$: 1 for *all*, 3/4 for *most*, 1/2 for *some*, 1/4 for *hardly any*, and 0 for *none*. To avoid confusion with the recall used in information retrieval, we call our metric weighted retention, $Retention_w$, and define it as follows:

$$\frac{(\text{Number of MUs marked}) \bullet C}{\text{Total number of MUs in the model summary}}$$
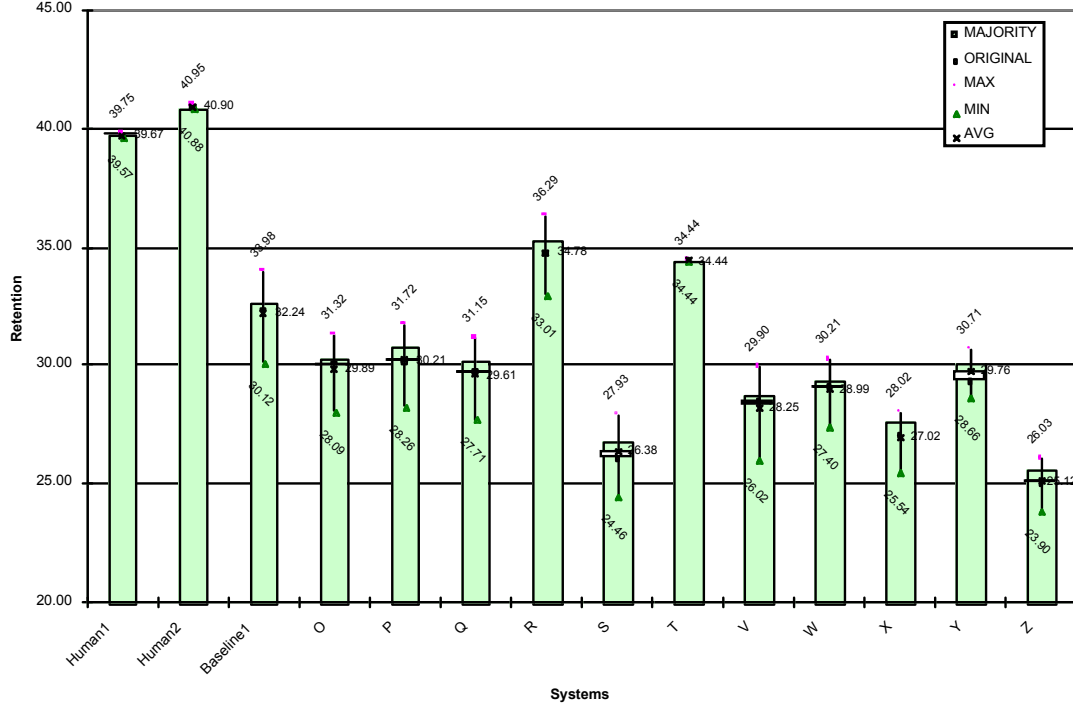
**Figure 2.** DUC 2001 single document retention score distribution.

If we ignore *C* (set it to 1), we obtain an unweighted retention, *Retention₁*. We used $Retention_1$ in our evaluation to illustrate that relative system performance (i.e., system ranking) changes when different evaluation metrics are chosen. Therefore, it is important to have common and agreed upon metrics to facilitate large scale evaluation efforts.

## 4    Instability of Manual Judgments

In the human evaluation protocol described in Section 2, nothing prevents an assessor from assigning different coverage scores to the same system units produced by different systems against the same model unit. (Since most systems produce extracts, the same sentence may appear in many summaries, especially for single-document summaries.)   Analyzing the DUC-2001 results, we found the following:

- Single document task

  o A total of 5,921 judgments

  o Among them, 1,076 (18%) contain multiple judgments for the same units

  o 143 (2.4%) of them have three different coverage scores

- Multi-document task

  o A total of 6,963 judgments

  o Among them 528 (7.6%) contain multiple judgments

  o 27 (0.4%) of them have three different coverage scores

Intuitively this is disturbing; the same phrase compared to the same model unit should always have the same score regardless of which system produced it.  The large percentage of multiple judgments found in the single document evaluation are test-retest errors that need to be addressed in computing performance metrics.

Figure 2 and Figure 3 show the retention scores for systems participating in the single- and multi-document tasks respectively.  The error bars are bounded at the top by choosing the maximum coverage score (MAX) assigned by an assessor in the case of multiple judgment scores and at the bottom by taking the minimum assignment (MIN).  We also compute system
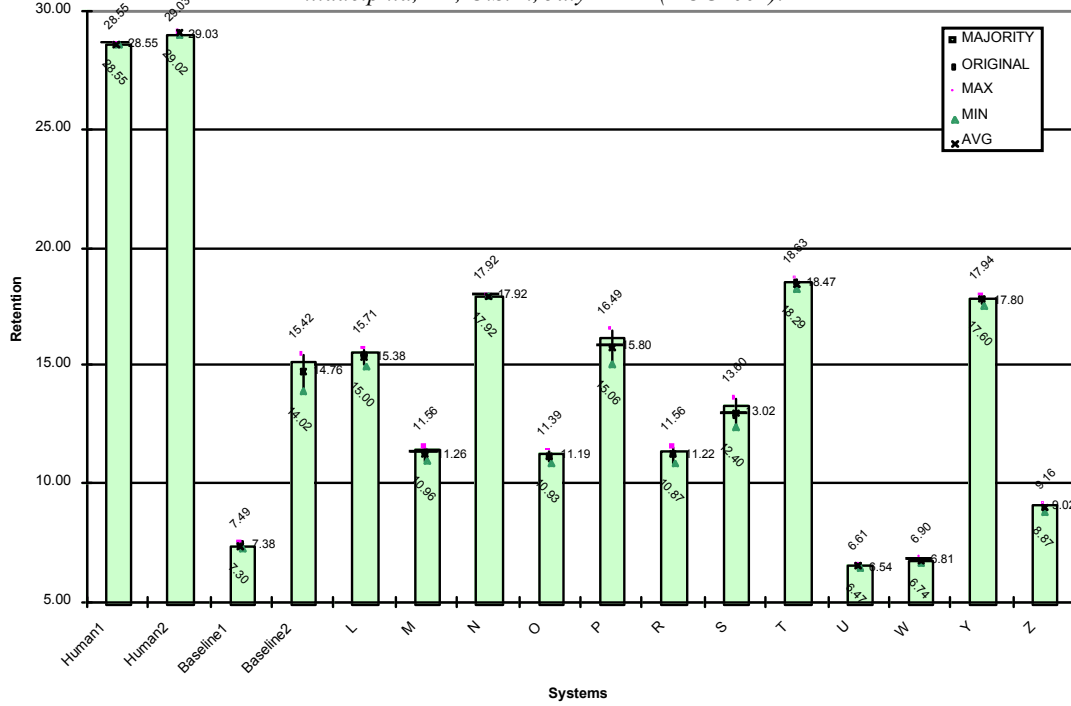
**Figure 3.** DUC 2001 multi-document retention score distribution.

retentions using the majority (MAJORITY) and average (AVG) of assigned coverage scores. The original (ORIGINAL) does not consider the instability in the data.

Analyzing all systems' results, we made the following observations.

(1) Inter-human agreement is low in the single-document task (~40%) and even lower in multi-documents task (~29%). This indicates that using a single model as reference summary is not adequate.

(2) Despite the low inter-human agreement, human summaries are still much better than the best performing systems.

(3) The relative performance (rankings) of systems changes when the instability of human judgment is considered. However, the rerankings remain local; systems remain within performance groups. For example, we have the following groups in the multi-document summarization task (Figure 3, considering 0.5% error):

    a. {Human1, Human2}
    b. {N, T, Y}
    c. {Baseline2, L, P}
    d. {S}
    e. {M, O, R}
    f. {Z}
    g. {Baseline1, U, W}

The existence of stable performance regions is encouraging. Still, given the large error bars, one can produce 162 different rankings of these 16 systems. Groups are less obvious in the single document summarization task due to close performance among systems.

Table 1 shows relative performance between systems *x* and *y* in the single document

|     | H1 | H2 | B1 | O | P | Q | R | S | T | V | W | X | Y | Z |
|-----|----|----|----|---|---|---|---|---|---|---|---|---|---|---|
| **H1** | = | - | + | + | + | + | + | + | + | + | + | + | + | + |
| **H2** | + | = | + | + | + | + | + | + | + | + | + | + | + | + |
| **B1** | - | - | = | ~ | ~ | ~ | ~ | + | - | + | ~ | + | ~ | + |
| **O** | - | - | ~ | = | ~ | ~ | - | + | - | ~ | ~ | + | ~ | + |
| **P** | - | - | ~ | ~ | = | ~ | - | + | - | ~ | ~ | + | ~ | + |
| **Q** | - | - | ~ | ~ | ~ | = | - | ~ | - | ~ | ~ | ~ | ~ | + |
| **R** | - | - | ~ | + | + | + | = | + | ~ | + | + | + | + | + |
| **S** | - | - | - | - | - | ~ | - | = | - | ~ | ~ | ~ | - | ~ |
| **T** | - | - | + | + | + | + | ~ | + | = | + | + | + | + | + |
| **V** | - | - | - | ~ | ~ | ~ | - | ~ | - | = | ~ | ~ | ~ | ~ |
| **W** | - | - | ~ | ~ | ~ | ~ | - | ~ | - | ~ | = | ~ | ~ | + |
| **X** | - | - | - | - | - | ~ | - | ~ | - | ~ | ~ | = | - | ~ |
| **Y** | - | - | ~ | ~ | ~ | ~ | - | + | - | ~ | ~ | + | = | + |
| **Z** | - | - | - | - | - | - | - | ~ | - | ~ | - | ~ | - | = |

**Table 1.** Pairwise relative system performance (single document summarization task).

|    | H1 | H2 | B1 | B2 | L | M | N | O | P | R | S | T | U | W | Y | Z |
|----|----|----|----|----|---|---|---|---|---|---|---|---|---|---|---|---|
| H1 | =  | -  | +  | +  | + | + | + | + | + | + | + | + | + | + | + | + |
| H2 | +  | =  | +  | +  | + | + | + | + | + | + | + | + | + | + | + | + |
| B1 | -  | -  | =  | -  | - | - | - | - | - | - | - | - | + | + | - | - |
| B2 | -  | -  | +  | =  | ~ | + | - | + | ~ | + | + | - | + | + | - | + |
| L  | -  | -  | +  | ~  | = | + | - | + | ~ | + | + | - | + | + | - | + |
| M  | -  | -  | +  | -  | - | = | - | ~ | - | ~ | - | - | + | + | - | + |
| N  | -  | -  | +  | +  | + | + | = | + | + | + | + | - | + | + | ~ | + |
| O  | -  | -  | +  | -  | - | ~ | - | = | - | ~ | - | - | + | + | - | + |
| P  | -  | -  | -  | ~  | ~ | + | - | + | = | + | + | - | + | + | - | + |
| R  | -  | -  | +  | -  | - | ~ | - | - | ~ | = | - | - | + | + | - | + |
| S  | -  | -  | +  | -  | - | + | - | + | - | + | = | - | + | + | - | + |
| T  | -  | -  | +  | +  | + | + | + | + | + | + | + | = | + | + | + | + |
| U  | -  | -  | -  | -  | - | - | - | - | - | - | - | - | = | - | - | - |
| W  | -  | -  | -  | -  | - | - | - | - | - | - | - | - | + | = | - | - |
| Y  | -  | -  | +  | +  | + | + | ~ | + | + | + | + | - | + | + | = | + |
| Z  | -  | -  | +  | -  | - | - | - | - | - | - | - | - | + | + | - | = |

**Table 2.** Pairwise relative system performance (multi-document summarization task).

summarization task. A '+' indicates the minimum retention score of *x* (row) is higher than the maximum retention score of *y* (column), a '-' indicates the maximum retention score of *x* is lower than the minimum retention score of *y*, and a '~' means *x* and *y* are indistinguishable. Table 2 shows relative system performance in the multi-document summarization task.

Despite the instability of the manual evaluation, we discuss automatic summary evaluation in an attempt to approximate the human evaluation results in the next section.

# 5   Automatic Summary Evaluation

Inspired by recent progress in automatic evaluation of machine translation (BLEU; Papineni et al. 2001), we would like to apply the same idea in the evaluation of summaries. Following BLEU, we used the automatically computed accumulative n-gram matching scores (NAMS) between a model unit (MU) and a system summary (S)[4] as performance indicator, considering multi-document summaries. Only content words were used in forming n-grams. NAMS is defined as follows:

$$a_1 \cdot NAM_1 + a_2 \cdot NAM_2 + a_3 \cdot NAM_3 + a_4 \cdot NAM_4$$

$NAM_n$ is n-gram hit ratio defined as:

$$\frac{\text{\# of matched n-grams between MU and S}}{\text{total \# of n-grams in MU}}$$

We tested three different configurations of $a_i$:

C1: $a_1 = 1$ and $a_2 = a_3 = a_4 = 0$;

C2: $a_1 = 1/3$, $a_2 = 2/3$, and $a_3 = a_4 = 0$;

C3: $a_1 = 1/6$, $a_2 = 2/6$, $a_3 = 3/6$, and $a_4 = 0$;

C1 is simply unigram matching. C2 and C3 give more credit to longer n-gram matches. To examine the effect of stemmers in helping the n-gram matching, we also tested all configurations with two different stemmers (Lovin's and Porter's). Figure 4 shows the results with and without using stemmers and their Spearman rank-order correlation coefficients (rho) compared against the original retention ranking from Figure 4. X-*n*G is configuration *n* without using any stemmer, L-*n*G with the Lovin stemmer, and P-*n*G with the Porter stemmer.

The results in Figure 4 indicate that unigram matching provides a good approximation, but the best correlation is achieved using C2 with the Porter stemmer. Using stemmers did improve correlation. Notice that rank inversion remains within the performance groups identified in Section 4. For example, the retention ranking of Baseline1, U, and W is 14, 16, and 15 respectively. The P-2G ranking of these three systems is 15, 14, and 16. The only system crossing performance groups is Y. Y should be grouped with N and T but the automatic evaluations place it lower, in the group with Baseline2, L, and P. The primary reason for Y's behavior may be that its summaries consist mainly of headlines, whose abbreviated style differs from the language models derived from normal newspaper text.

For comparison, we also ran IBM's BLEU evaluation script[5] over the same model and system summary set. The Spearman rank-order correlation coefficient ($\rho$) for the single document task is 0.66 using one reference summary and 0.82 using three reference summaries; while Spearman $\rho$ for the multi-document task is 0.67 using one reference and 0.70 using three.

# 6   Conclusions

We described manual and automatic evaluation of single and multi-document summarization in DUC-2001. We showed the instability of

---

[4] The whole system summary was used to compute NAMS against a model unit.

[5] We thank Kishore Papineni for sending us BLEU 1.0.

| SYSCODE | Original Retention ranking | No stemmer | | | Lovin stemmer | | | Porter stemmer | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | X-1G (unigram) | X-2G | X-3G | L-1G (unigram) | L-2G | L-3G | P-1G (unigram) | P-2G | P-3G |
| Human1 | 1 | 2 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| Human2 | 2 | 1 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 |
| Baseline1 | 14 | 15 | 15 | 15 | 16 | 15 | 14 | 16 | 15 | 14 |
| Baseline2 | 8 | 8 | 7 | 6 | 8 | 8 | 6 | 8 | 8 | 6 |
| L | 7 | 7 | 6 | 7 | 7 | 7 | 7 | 7 | 7 | 7 |
| M | 10 | 10 | 10 | 10 | 10 | 11 | 11 | 9 | 10 | 11 |
| N | 4 | 4 | 4 | 4 | 4 | 4 | 4 | 4 | 4 | 4 |
| O | 11 | 12 | 12 | 12 | 12 | 12 | 12 | 11 | 12 | 12 |
| P | 6 | 5 | 5 | 5 | 5 | 5 | 5 | 5 | 5 | 5 |
| R | 11 | 11 | 11 | 11 | 11 | 10 | 10 | 12 | 11 | 10 |
| S | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 |
| T | 3 | 3 | 3 | 3 | 3 | 3 | 3 | 3 | 3 | 3 |
| U | 16 | 14 | 14 | 14 | 14 | 14 | 15 | 14 | 14 | 15 |
| W | 15 | 16 | 16 | 16 | 15 | 16 | 16 | 15 | 16 | 16 |
| Y | 5 | 6 | 8 | 8 | 6 | 6 | 8 | 6 | 6 | 8 |
| Z | 13 | 13 | 13 | 13 | 13 | 13 | 13 | 13 | 13 | 13 |
| Spearman ρ | 1.00000 | 0.98382 | 0.97206 | 0.96912 | 0.98382 | 0.98382 | 0.97206 | 0.98235 | 0.98676 | 0.97206 |

**Figure 4.** Manual and automatic ranking comparisons.

human evaluations and the need to consider this factor when comparing system performances. As we factored in the instability, systems tended to form separate performance groups. One should treat with caution any interpretation of performance figures that ignores this instability.

Automatic evaluation of summaries using accumulative n-gram matching scores seems promising. System rankings using NAMS and retention ranking had a Spearman rank-order correlation coefficient above 97%. Using stemmers improved the correlation. However, satisfactory correlation is still elusive. The main problem we ascribe to automated summary evaluation is the large expressive range of English since human summarizers tend to create fresh text. No n-gram matching evaluation procedure can overcome the paraphrase or synonym problem unless (many) model summaries are available.

We conclude the following:

(1) We need more than one model summary although we cannot estimate how many model summaries are required to achieve reliable automated summary evaluation.

(2) We need more than one evaluation for each summary against each model summary.

(3) We need to ensure a single rating for each system unit.

## References

DUC. 2001. The Document Understanding Conference 2001. http://www-nlpir.nist.gov/ projects/duc/2001.html.

Lin, C.-Y. 2001. *Summary Evaluation Environment*. http://www.isi.edu/~cyl/SEE.

Mani, I., D. House, G. Klein, L. Hirschman, L. Obrst, T. Firmin, M. Chrzanowski, and B. Sundheim. 1998. *The TIPSTER SUMMAC Text Summarization Evaluation: Final Report*. MITRE Corp. Tech. Report.

Papineni K., S. Roukos, T. Ward, W.-J. Zhu. 2001. *BLEU: a Method for Automatic Evaluation of Machine Translation*. IBM Research Report RC22176(W0109-022).