

**Stanford University Computer Science Department  
CS347 Spring 2001 Mid-term (Total of 30 points.)**

This exam is open book, open notes. You have 70 minutes.

Print your name: \_\_\_\_\_

The Honor Code is an undertaking of the students, individually and collectively:

1. that they will not give or receive aid in examinations; that they will not give or receive unpermitted aid in class work, in the preparation of reports, or in any other work that is to be used by the instructor as the basis of grading;
2. that they will do their share and take an active part in seeing to it that others as well as themselves uphold the spirit and letter of the Honor Code.

The faculty on its part manifests its confidence in the honor of its students by refraining from proctoring examinations and from taking unusual and unreasonable precautions to prevent the forms of dishonesty mentioned above. The faculty will also avoid, as far as practicable, academic procedures that create temptations to violate the Honor Code.

While the faculty alone has the right and obligation to set academic requirements, the students and faculty will work together to establish optimal conditions for honorable academic work.

I acknowledge and accept the Honor Code.

Signed: \_\_\_\_\_

Problem	Points	Maximum
1		4
2		5
3		5
4		5
5		6
6		5
<b>Total</b>		



**Question 3:**

(5 points)

Recall the estimate of the total size of the postings entries (45Mbytes using 2 codes) from Lecture 1 using Zipf's law. Using the same parameters (1 million documents, 500,000 terms), re-compute this estimate if we were to omit from indexing the 1% of the most frequently occurring terms.

**Question 4:**

(5 points)

Mark each of the following assertions as True or False.

	Assertion	T/F
(i)	The optimal order for query processing in an <i>AND</i> query is always realized by starting with the term occurring in the fewest documents.	
(ii)	The 2 code for 17 is 111000001.	
(iii)	The base of the logarithm used in the $\text{tf-idf}$ formula makes no difference to the cosine distance between two documents (provided they both use the same base).	
(iv)	If we were to take a document and double its length by repeating every occurrence of every word, then the normalized $\text{tf-idf}$ values for all terms in this document remain unchanged.	
(v)	The optimal order for query processing in an <i>AND</i> query is always realized by starting with the term occurring in the fewest documents.	