

15-440 Distributed Systems Final Exam

Name:

Andrew: ID

December 12, 2011

- Please write your name and Andrew ID above before starting this exam.
- This exam has 14 pages, including this title page. Please confirm that all pages are present.
- This exam has a total of 80 points.

Question	Points	Score
1	8	
2	3	
3	6	
4	12	
5	10	
6	12	
7	13	
8	6	
9	10	
Total:	80	

Short Answers

1. (8 points) Circle True or False as appropriate. If you don't know the answer, come back at the end and GUESS; a blank answer is the same as a wrong answer, so guessing can only help.

True False DNS is delegated hierarchically by having, e.g., the root nodes tell resolvers which servers to query for “.com”, and so on, until the client's query can be answered.

True False Domain Name System (DNS) resolvers use Paxos and invalidation messages to maintain the consistency of cached records.

True False In Tor's “onion routing”, a compromised node in the middle of the network can identify the client and destination web site of the communication.

True False If the server a client communicates with does not support any form of encryption (e.g., https), the client should use Tor so that nobody can overhear its traffic.

True False Some Content Delivery Networks (CDNs) use DNS responses to direct clients to the closest CDN cache.

True False A major challenge for Tor is finding nodes to act as exit points from the Tor network, because these nodes may appear to be performing illegal activities on the Internet.

True False The goal of paravirtualization (e.g., Xen) is to make the guest operating system completely unaware that it is running on a virtual machine.

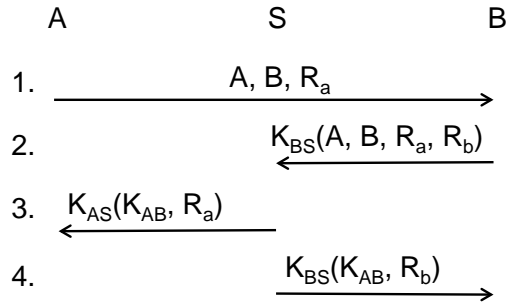
True False The Map step of MapReduce provides a way to store and retrieve items according to their keys.

2. (3 points) Cloud computing services like Amazon's EC2 assign users virtual machines (VMs) instead of allocating physical machines directly. Doing so provides at least three major benefits to Amazon. Explain what these three benefits are, giving a brief motivation for each one.

Security Protocols

4. (12 points) The following figure illustrates a protocol between two agents A and B, and a server S. The intent is to get the server to generate a shared session key K_{AB} , and to also enable A and B to make sure they are communicating with each other.

Each agent shares a private key with the server S: A has K_{AS} , and B has K_{BS} . Values R_a and R_b are randomly generated “nonces” (number-used-once). The notation $K(M_1, M_2, \dots, M_n)$ means to generate a message containing an encrypted version of the sequence M_1, M_2, \dots, M_n using key K .



Assume the following:

- Initially, only A and S know K_{AS} , and only B and S know K_{BS} .
- The true server S is not malicious, but there could be an imposter S' trying to pose as S.
- There could be an imposter A' or B' trying to pose as A or B.
- Imposters can intercept any traffic, replay old messages, or inject new ones.
- The encryption is secure, and the encrypted form of the sequence does not reveal any information about the encrypted form of the individual elements. For example, knowing $K(M_1, M_2)$ does not reveal any information about $K(M_1)$ or $K(M_2)$.
- We will say that a message is *freshly generated* if it must have been created some time after the initial message in the protocol.

For each of the following statements, state whether it is true or false. Give a brief (one or two sentences should suffice) justification for your answer.

- (a) S can be certain that message 2 was freshly generated by B.

- (b) A can be certain that message 3 was freshly generated by S.

- (c) B can be certain that message 4 was freshly generated by S.

- (d) Upon completion of the protocol, A can be certain that it has established a session with B.

- (e) Upon completion of the protocol, B can be certain that it has established a session with A.

- (f) Upon completion of the protocol, no one other than A, B, or S can know the value of K_{AB} .

Distributed Replication

5. (a) (2 points) Consider an implementation of the Paxos algorithm where a leader waits for less than a majority of acceptors to answer OK to a Prepare or an Accept message before it proceeds to executing the next steps. Explain how such an implementation can break the Paxos guarantees.

(b) (2 points) Explain why Paxos cannot tolerate f failures with less than $2f + 1$ nodes.

(c) (6 points) One common assumption in Paxos is that different leaders use different proposal numbers. A colleague shows you an implementation of Paxos where:

1. Different leaders can use the same proposal numbers, and
2. An acceptor rejects a Prepare or Accept message only if the proposal number in the message is strictly smaller than the largest proposal number that acceptor has seen (i.e., exactly as presented in class and in the course notes).

Is this Paxos implementation correct? Find a counterexample or explain briefly why it is correct.

Peer to Peer

6. (a) (2 points) In a centralized p2p network (such as the old Napster), how many indices must be searched if a client wants to locate a particular file?
- (b) (2 points) Name one major disadvantage of the centralized p2p system (from a distributed principles point of view)
- (c) (2 points) Query flooding is an alternative design that solves some of the problems of centralized p2p and eliminates the central server. However, it changes the mechanics of peer interactions significantly. Explain (1 sentence each) how a newly-joining node publishes the files they wish to make available, in...

A centralized p2p network:

A query flooding p2p network:

- (d) (2 points) One popular improvement upon query flooding is to move to a “supernode” flooding architecture. Using N as the number of nodes in the network and, S as the number of supernodes ($S \ll N$), explain the benefit of moving to this supernode architecture.

(e) (2 points) What is a common mechanism used to limit the propagation of queries in a flooding network?

(f) (2 points) List one typical criterion for selecting a node to be promoted to a supernode. Explain in one sentence why such a choice would improve network stability.

Consistent Hashing

7. (a) (2 points) In a distributed hash table (DHT) using a *ring mapping*, explain how a key is mapped to a specific node?
- (b) (2 points) In the most basic form of DHT systems, nodes simply track their predecessor and successor. Very briefly state what problem this leads to as the network grows larger.
- (c) (2 points) The Chord scheme overcomes this problem by having each node maintain a “finger table”. Given starting node j and destination node k , where $k > j$, how much is the distance between the two nodes reduced with each hop?
- (d) (2 points) In a Chord structured DHT with N nodes, how many hops would a lookup operation require?
- (e) (2 points) Assume a Chord network in which node q is the successor of node p . During operation, node p discovers that its successor link is no longer consistent because q has updated its predecessor link. What does this change imply must have happened in the network?

- (f) (3 points) An optimization to Chord involves storing several nodes for each entry in the finger table instead of just one. Explain an important benefit this optimization confers in a globally distributed DHT.

Byzantine Fault Tolerance

8. (6 points) In distributed systems where messages are asynchronous and failures can be Byzantine, we have to use at least $n = 3f + 1$ replicas in total to tolerate f faulty replicas. Show that this bound is tight, i.e., that $n \geq 3f + 1$ *must* hold in order for the system to work properly.

To approach this problem, consider this scenario: a client sends the same command to each of n servers and then waits for the servers to execute the command and send back the result. Ideally, the client should receive n matching results, but remember that f servers may be malicious. Additionally, messages are asynchronous, so they can be delayed for an indefinite amount of time. Show that $n \geq 3f + 1$ must hold for the client to always be able to identify the correct result.

MapReduce and GFS

9. (10 points) MapReduce has proved an extremely popular framework for distributed computation on large clusters, because it masks many of the painful parts of ganging together thousands of nodes to accomplish a task.
- (a) In general high-performance computing (HPC), programmed using message passing, what is the technique used to provide fault tolerance? (Very short answer)

 - (b) Identify two important limitations that MapReduce places upon the Map functions so that the framework can hide failures from the programmer.

 - (c) What advantage does this give MapReduce over MPI for failure handling and recovery? (There are several—list the one(s) you think is most important)

 - (d) How does MapReduce handle “straggler” tasks that take longer than all of the others (e.g., perhaps they’re running on a slower machine or one with buggy hardware)?

 - (e) MapReduce and the Google File System (GFS) were designed to work well together. What important optimization in MapReduce is enabled by having GFS expose block replica locations via an API?