

MIT OpenCourseWare  
<http://ocw.mit.edu>

6.033 Computer System Engineering  
Spring 2009

For information about citing these materials or our Terms of Use, visit: <http://ocw.mit.edu/terms>.

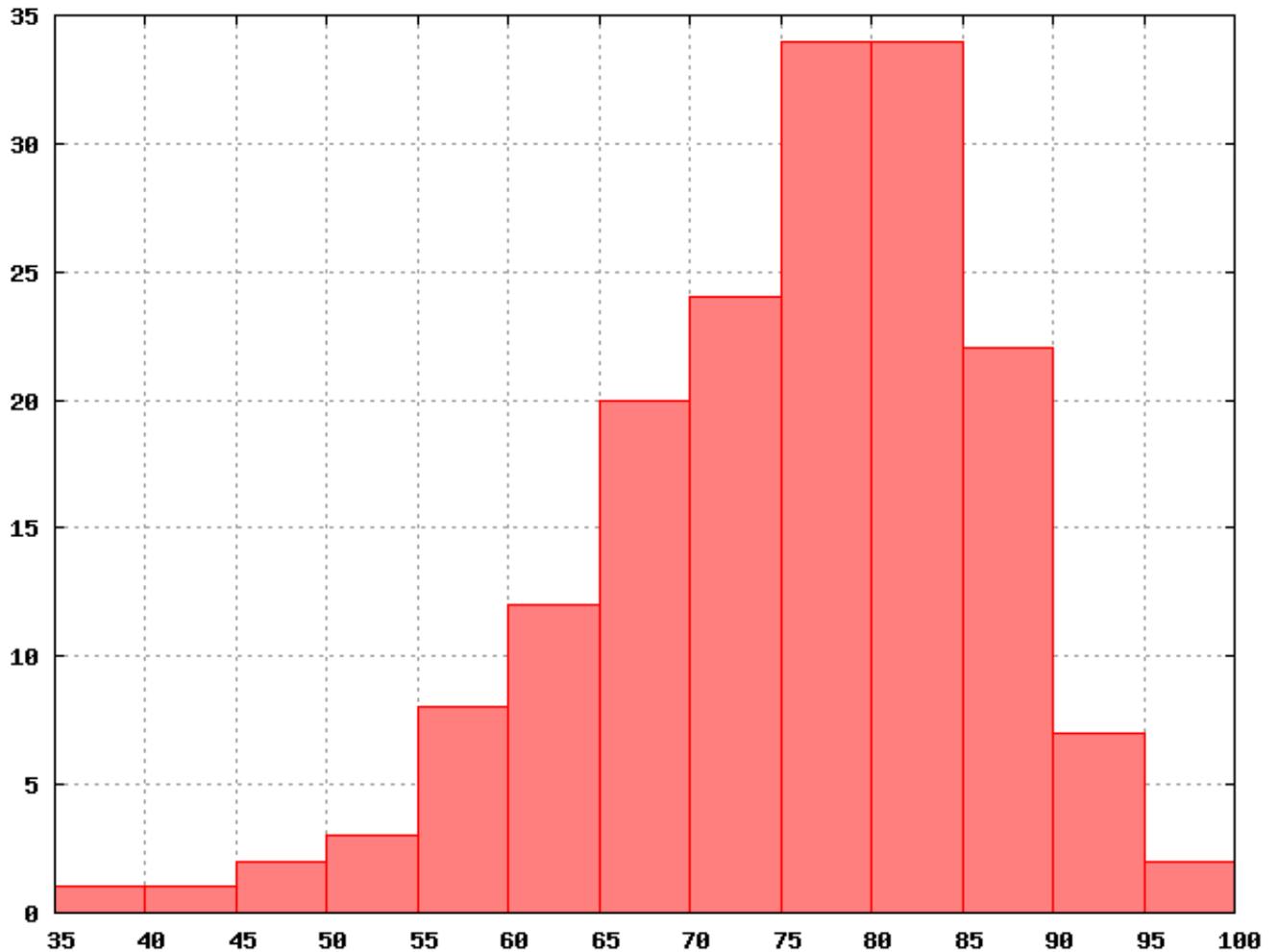


*Department of Electrical Engineering and Computer Science*

MASSACHUSETTS INSTITUTE OF TECHNOLOGY

**6.033 Computer Systems Engineering: Spring 2006**

## Quiz II



$N = 170$ , MEAN =  $75.1/96 = 78\%$ , MEDIAN =  $76/96 = 79\%$ , STDDEV = 10.7

## I Reading Questions

**1. [6 points]:** Ethernet provides a medium access protocol (MAC) that allows multiple machines to communicate using a shared medium (Reading #9. Robert M. Metcalfe and David R. Boggs. Ethernet: Distributed packet switching for local computer networks). TDMA is a different MAC protocol that has been presented in lectures. Say that there are  $n$  machines connected to the same bus. A TDMA approach would divide time into frames. Each frame has  $n$  slots. Each slot is sufficient to send one packet. Each one of the  $n$  machines is assigned a number between 1 and  $n$ . If machine  $i$  has a packet to send, it waits for slot  $i$  and sends its packet.

**(Circle True or False for each choice.)**

*This question was graded +2 for a correct answer, -1 for an incorrect answer, and 0 for no answer.*

**A. True / False** Ethernet achieves a higher data rate than TDMA when a small number of machines have packets to send (i.e., much smaller than  $n$ ), whereas TDMA is more efficient if most of the time all  $n$  machines have packets to send.

*TRUE. With TDMA, each machine gets  $1/n$  of the total available bandwidth — no more, no less. With Ethernet, any machine can send successfully as long as there are no collisions. When a small number of machines try to send, they will get most of the bandwidth, not just  $1/n$  of the bandwidth per machine.*

**B. True / False** Ethernet provides a guaranteed minimum throughput for each of the  $n$  machines

*FALSE. Because of the possibility of collisions, there is no absolute guarantee that any machine will be able to send any packet successfully.*

**C. True / False** Ethernet guarantees no collision because each machine senses the carrier before sending and does not send if another machine is currently using the medium.

*FALSE. If two machines start sending at the same time, neither will sense the other because the signals have not had time to propagate from one machine to the other. In this case there will be a collision.*

**2. [8 points]:** Assume that an NFS (described in appendix 4.B) server contains a file `/a/b` and that an NFS client mounts the NFS server's root directory in the location `/x`, so that the client can now name the file as `/x/a/b`. Further assume that this is the only client and that the following code executes on the client:

```
chdir("/x/a"); // change to directory /x/a
rm("b");
```

The REMOVE message from the client to the server gets through and the server removes the file. Unfortunately, the response from the server to the client is lost so the client resends the message to remove the (now non-existent) file. The server receives the resent message. What happens next depends on the server implementation. Which of the following are correct statements?

**(Circle True or False for each choice.)**

*This question was graded +2 for a correct answer, -1 for an incorrect answer, and 0 for no answer.*

Name:

- A. True / False** If the server maintains an in-memory reply cache in which it records all operations it previously executed, and there are no server failures, the server will return “OK”.
- TRUE. When the server will get the resent message, it will look in its reply cache, discover that it has already processed the message, and return the same response to the client. Note that the question clearly indicated that the server recorded all previously executed operations. Note that in practice, no server can ever afford to guarantee to record all operations in its memory reply cache because that would force the server to keep a potentially unbounded amount of state around forever. One way around this problem is to keep a fixed amount of space for the cache and discard results of old operations when space gets tight.*
- B. True / False** If the server maintains an in-memory reply cache but the server has failed, restarted, and its reply cache is empty, both of the following responses are possible: the server may return “file not found” or “OK”.
- FALSE. If the server processes the resent message with an empty reply cache, it will always try to execute the request in the message. In this case it will always return “file not found” because it removed the file when it processed the first request. The file is therefore no longer present.*
- C. True / False** If the server is stateless, it will return “file not found”.
- TRUE. A stateless server will not maintain a reply cache, so it will try to execute whatever request comes in. In particular, it will try to remove the file, find that it is not present, and return “file not found.”*
- D. True / False** Because REMOVE is an idempotent operation, any server implementation will return “OK”.
- FALSE. As the previous discussion indicates, REMOVE is not an idempotent operation in NFS.*

**3. [6 points]:** Ben Bitdiddle is performing an important *deterministic computation* that will take roughly 2 years to complete. In order to add some fault tolerance, he buys 30 identical computers and runs the same program with the same inputs on all of them. When they return their results after computing for 2 years, the voter selects the majority answer (voting is described in Chapter 8). Which of the following failures can this scheme tolerate, assuming the voter works correctly?

**(Circle True or False for each choice.)**

*This question was graded +3 for a correct answer, -1 for an incorrect answer, and 0 for no answer.*

**A. True / False** The software carrying out the deterministic computation might have a bug in it, which causes it to compute the wrong answer for certain inputs.

*FALSE. If the software has a deterministic error, all versions will compute the same incorrect result.*

**B. True / False** Cosmic rays might corrupt data stored in memory at one or two computers, causing them to return incorrect results.

*TRUE. Even though the computers with corrupted memory may compute an incorrect result, the majority of the computers will compute the correct result. The voting procedure will ensure that Ben gets the correct result.*

Consider the following protocol: The sender assigns each packet a sequence number. The receiver sends cumulative acks, i.e, whenever it receives a packet, the receiver sends an ack for the next expected sequence number. For example, when it receives packets 1, 2, 4, the receiver sends acks 2, 3, 3 respectively. The sender uses a sliding window protocol. The size of the sliding window is *fixed* to 4 packets (because the receiver cannot buffer more than 4 packets at once). The network is given below. The data rate between the sender's computer and router is 200 packets/second, and the rate between router and receiver's computer is 100 packets/second. The one-way propagation delay between sender and receiver is 0.05 second (i.e., the round trip propagation delay is 0.1 second).

Sender ----- Router ----- Receiver  
           200 packets/second                  100 packets/second

**4. [2 points]:** Given the above protocol and network, what is the maximum number of packets the sender can send in a burst without waiting for acks?

*The correct answer is 4 packets because in a window protocol the sender can send a whole window without waiting for acks, but cannot send more than a window.*

**5. [4 points]:** Given the above protocol and network, what is the maximum send rate measured in packets per second?

*The correct answer is 40 packets/sec, which is the window size divided by the round trip time. We have given partial credit to people who said 100 packets per second, which is the bottleneck capacity along the path. We have also given partial credit to people who divided by the one-way delay instead of the RTT.*

**6. [6 points]:** Assume that at time  $t=0$ , the sender sent packets 1,2,3,4. By time  $t=0.15$ s, the sender has received acks 2,2,2 and nothing more. Which of the following statements are correct?

**(Circle True or False for each choice.)**

*This question was graded +2 for a correct answer, with no penalty for incorrect answers.*

**A. True / False** The receiver has not received packet 4

*FALSE. The fact that the sender received acks 2,2,2 means that the receiver received 3 distinct packets and is still missing packet 2. Thus, the receiver has received packet 4.*

**B. True / False** The sender is sure that packet 1 has been received.

*TRUE. The fact that the sender received acks 2,2,2 means that the receiver received the first packet.*

**C. True / False** By  $t=0.15$ , the sliding window has advanced by 3 packets

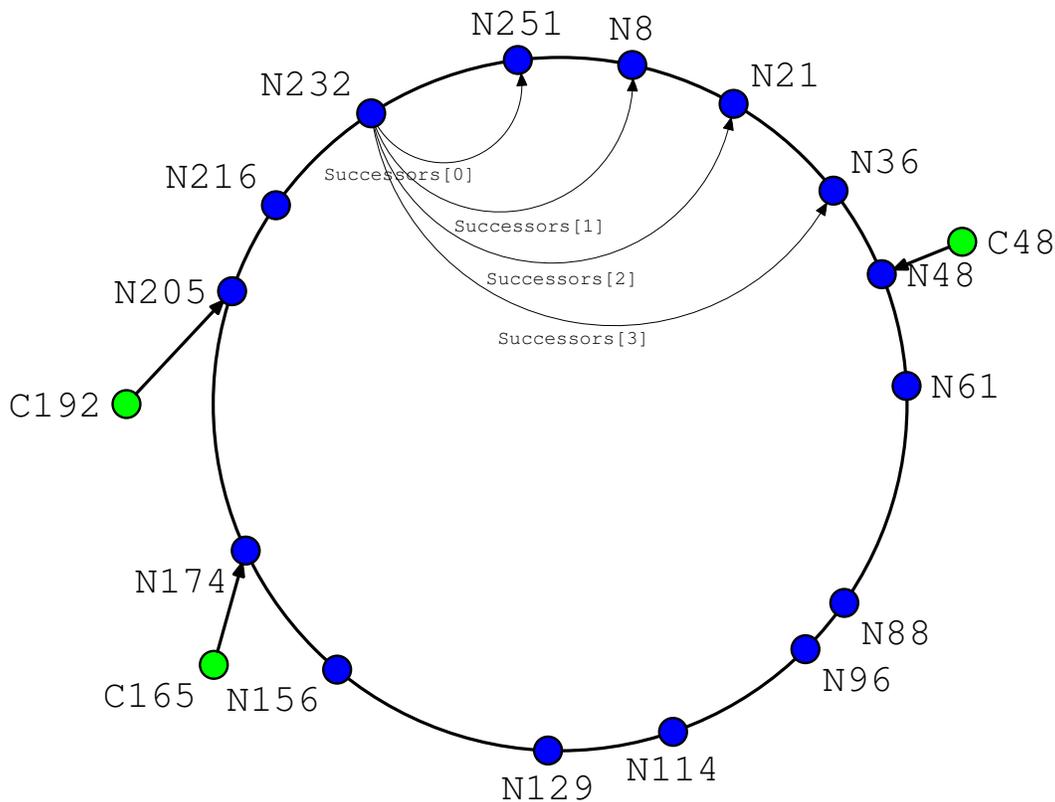
*FALSE. The window cannot advance beyond packet two until it is ack-ed.*

Name:

## II SURETHING

Alyssa P. Hacker decides to offer her own content delivery system, SURETHING. A SURETHING system contains 1000 computers that communicate over the Internet. Each computer has a unique ID, and they are thought of as (logically) being organized in a ring (see Figure 1). Each computer has *successors* as shown in the figure. The ring “wraps-around”: the immediate successor of the computer with the highest ID (computer 251 in the figure) is the computer with the lowest ID (computer 8).

Each content item also has a unique ID,  $c$ , and the content will be stored at  $c$ 's *immediate successor*: the first computer in the ring whose ID exceeds  $c$  (see Figure 1).



**Figure 1:** Arrangement of computers in a ring. Computer #232's pointers to its 4 successors are shown. The content item  $c$  with id #C192 will be stored at its immediate successor in the ring, computer #N205; item number #C48 will be stored at its immediate successor, computer #N48; and item number #C165 will be stored at its immediate successor, computer #N174.

Alyssa designs the system using two layers: a forwarding and routing layer (to find the IP address of the computer that stores the content) and a content layer (to store or retrieve the content).

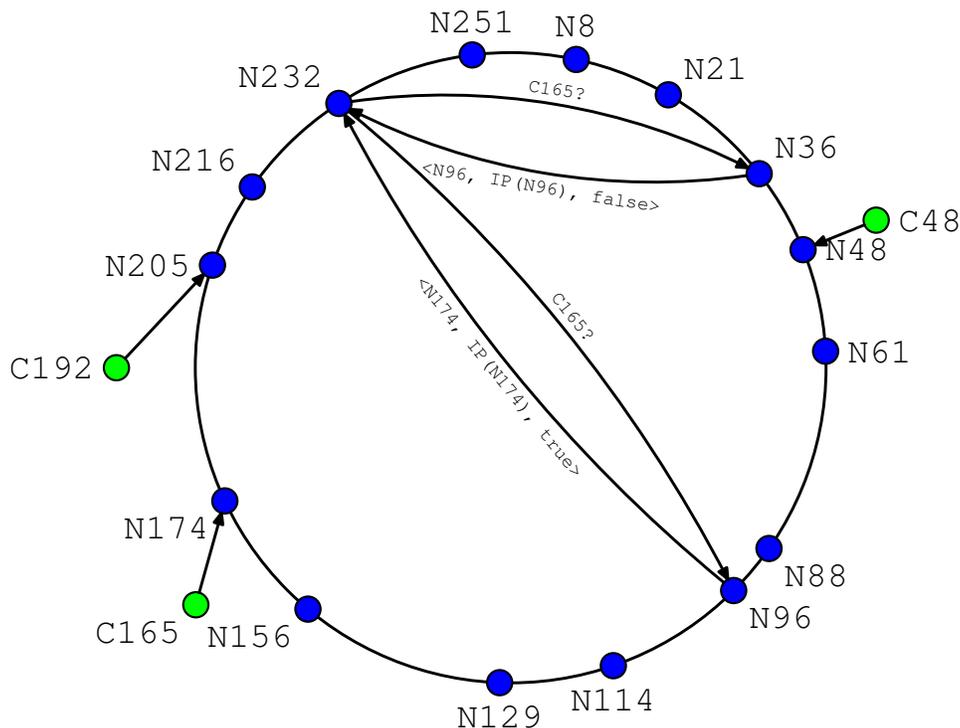
Name:

## II.1 Building a Forwarding and Routing Layer

Initially Alyssa decides that the routing step will work as follows: Each computer has a local table, *Successors*, that contains the ID and IP address of its 4 successors (the 4 computers whose IDs follow this computer's ID in the ring); the entries are ordered as they appear in the ring. These tables are set up when the system is initialized; you can assume they contain the right information. Unless a question states otherwise, you may assume that computers don't fail.

The forwarding and routing layer of each node provides a procedure GETLOC that can be called by the content layer to find the IP address of the immediate successor of some content item  $c$ . This procedure checks its local *Successors* table to see if it contains the immediate successor of the requested content; if not, it makes a remote procedure call to the GETLOCLocal procedure on the *most distant* successor in its *Successors* table. That computer returns the immediate successor of  $c$  if it is known locally in its *Successors* table; otherwise that node returns its *most distant* successor, and the originating computer continues the search there, iterating in this way until it locates  $c$ 's immediate successor.

For example, if computer N232 is looking for the immediate successor of  $c = C165$  in the system shown in Figure 1, it will first look in its local table; since this table doesn't contain the immediate successor of  $c$ , it will request information from computer N36. Computer N36 also doesn't have the immediate successor of  $C165$  in its local *Successors* table, and therefore it returns the IP address of computer N96. Computer N96 does have the immediate successor (computer N174) in its local *Successors* table and it returns this information. This sequence of RPCs is shown in Figure 2.



**Figure 2:** Sequence of RPCs and replies required for computer N232 to find the immediate successor of the content item with ID C165.

Name:

The code for GETLOC is given in Figure 3. This code correctly implements the approach described in the preceding paragraph. The code uses a helper procedure, LOOKUP, to search the *Successors* table; LOOKUP returns the first computer in the table that is a successor of the item of interest, or the last entry in the table, if the table doesn't contain a successor.

The code makes use of an RPC mechanism, invoked by DOCALL, that takes care of sending the message (including re-sending it if necessary), calling the requested procedure (GETLOCLocal) at the receiving end, getting the result from this call, and sending it back to the caller. The RPC mechanism runs at the application layer and uses UDP as the transport layer to actually do the communication.

```

structure pair { int id; IP loc; }
structure res { int id; IP loc; bool succ };
structure msg { int c; int id; IP loc; bool found; int resultCode}
define succsSz 4;
pair Successors[succsSz];

// the procedure called by the user
procedure GETLOC(int c)
  res ans ← LOOKUP(Successors, succsSz, c); // check local table
  if ans.found then return ⟨ans.id, ans.loc⟩; // answer found locally
  msg m, reply;
  m.c ← c;
  while TRUE do { // search iteratively
    reply ← DOCALL(ans.loc, GETLOCLocal, m);
    if reply.resultCode = OK then {
      if reply.found then return ⟨reply.id, reply.loc⟩; // remote computer found answer
      ans.loc ← reply.loc; // remote computer returned its largest successor; continue searching there
    } else // error handling (not shown)
  }
}

// the handler invoked by a RPC request when call is to GETLOCLocal
procedure GETLOCLocal(msg req)
  msg reply;
  res ans ← LOOKUP(Successors, succsSz, req.c); // check local table
  reply.resultCode ← OK; reply.loc ← ans.loc; reply.id ← ans.id; reply.found ← ans.found;
  return reply;

procedure LOOKUP(pair[] table, int sz, int c)
  if c = myID then return ⟨myID, myIP, true⟩; // check the local node
  for i ← 0 to sz - 1 do {
    if BETWEEN(c, myID, table[i].id) then return ⟨table[i].id, table[i].loc, true⟩;
    // between(c,x,y) returns true if c in (x, y) on the ring
  }
  return ⟨table[sz - 1].id, table[sz - 1].loc, false⟩; // return most distant successor

```

**Figure 3:** Alyssa's initial implementation for GETLOC

**Name:**

**7. [6 points]:** Consider the packet corresponding to an RPC message when it is traversing an Ethernet link. What headers will be added? Draw a diagram of the packet with these headers. You do not need to indicate the specific contents of each part, but should label each header with the layer it belongs to.

<i>ETHERNET (or LINK-LAYER) header</i>	<i>IP (or NETWORK-LAYER or ROUTING/ FORWARDING LAYER) header</i>	<i>UDP (or END-TO-END or APPLICATION LAYER) header</i>	<i>RPC header</i>	<b>RPC payload</b>
--	--	--	-------------------	------------------------

+2 points

+2 points

+2 points

We didn't deduct  
if the RPC header  
was omitted.

Wrong order: -2 points

"Content header": -1 point. These were dealt with case-by-case.

**8. [8 points]:** While testing SURETHING, Alyssa notices that when the Internet attempts to deliver the RPC packets, they don't always arrive at their destination. Which of the following reasons might prevent a packet from arriving at its destination?

**(Circle True or False for each choice.)**

*This question was graded +2 for a correct answer, -1 for an incorrect answer, and 0 for no answer.*

**A. True / False** A router discards the packet.

*TRUE. If an intermediate router discards a packet, the packet will not arrive at its destination.*

**B. True / False** The packet is corrupted in transit.

*TRUE. If a packet is corrupted in transit, the IP checksum may fail, causing an intermediate router to discard the packet.*

**C. True / False** The payload of the message contains the wrong *loc*

*FALSE. loc is part of the RPC payload, which is itself part of the IP (Internet Protocol) payload; Internet routers act on the IP header. Thus, the value of loc does not affect Internet forwarding.*

**D. True / False** The packet gets into a forwarding loop in the network.

*TRUE. Forwarding loops in the Internet are a fact of life. The purpose of the time-to-live field in IP packets is exactly to detect this case.*

For the next two questions, remember that computers don't fail and that all tables are initialized correctly.

**9. [6 points]:** Assume that  $c$  is an id whose immediate successor is not present in *Successors*, and  $n$  is the number of computers in the system. In the *best case*, how many remote lookups are needed before GETLOC( $c$ ) returns?

(Circle the BEST answer)

A. 0

B. 1

*This is the correct answer. In the best case, the first computer that GETLOC queries has in its Successors table the immediate successor of  $c$ , so only one remote lookup is required.*

C. 2

D.  $O(\log n)$

E.  $O(n)$

F.  $O(n^2)$

**10. [6 points]:** Assume that  $c$  is an id whose immediate successor is not present in *Successors*, and  $n$  is the number of computers in the system. In the *worst case*, how many remote lookups are needed before GETLOC( $c$ ) returns?

(Circle the BEST answer)

A. 0

B. 1

C. 2

D.  $O(\log n)$

E.  $O(n)$

*This is the correct answer. In the worst case, the computer executing GETLOC is "one past" the immediate successor of  $c$ . In that case, GETLOC must query every fourth computer in the system. Thus, the number of lookups needed is  $n/4$ , which is  $O(n)$ .*

F.  $O(n^2)$

## II.2 Building the Content Layer

Having built the forwarding and routing layer, Alyssa turns to building a content layer. At a high level, the system supports storing data that has an ID associated with it. Specifically, it supports two operations:

- $PUT(c, content)$  stores *content* in the system with ID *c*.
- $GET(c)$  returns the content that was stored with ID *c*.

You can assume that content IDs are integers that can be used as arguments to  $GETLOC$ . (In practice, this can be assured by using a hash function that maps human-readable names to integers.)

Alyssa implements the content layer by using the forwarding and routing layer to choose which computers to use to store the content. For reliability, she decides to store every piece of content on two computers: the two immediate successors of the content's ID. She modifies  $GETLOC$  to return both successors, calling the new version  $GETLOC2$ . For example, if  $GETLOC2$  is asked to find the successors of the content item with ID C165, it returns the IP addresses of computers N174 and N205 (given the information shown in Figure 1).

Once the correct computers are located using the forwarding and routing layer, Alyssa's implementation sends a  $REMOTEPUT$  RPC (which invokes the procedure  $REMOTEPUT$  on the remote computer) to each of these computers, requesting that the remote computer store the content in a file on its disk. To retrieve the content associated with a given ID, it sends a  $REMOTEGET$  RPC (which invokes the procedure  $REMOTEGET$  on the remote computer), requesting that the computer load the appropriate file from disk, if it exists, and return its contents.

The code for Alyssa's of  $GET$  and  $PUT$  is shown in Figure 4. It is not necessary to examine this code very closely, and you may assume it is correct. As before,  $DOCALL$  sends an RPC message, using retransmissions as necessary to ensure that the message is delivered.  $FETCH$  and  $STORE$  are an interface to the local file system; they maintain files so that previously stored content can be found, and return error codes when there are problems, e.g., if the requested content is not present.  $REMOTEGET$  and  $REMOTEPUT$  are't shown in the figure; you can assume they use  $FETCH$  and  $STORE$  to obtain or store the requested content.

```

structure res { int id; IP loc; bool succ };
structure cmsg { int c; data content; int resultCode; }

// called by a user PUT request
procedure PUT(int c, data content)
  cmsg m; m.c ← c; m.content ← content;
  res ⟨ans1, ans2⟩ ← GETLOC2(c)
  for ans in ⟨ans1, ans2⟩ do { // add to each of the two successors
    if (ans.id = myID) then STORE(c, content)
    else DOCALL(ans.loc, REMOTEPUT, m);
  }

// called by a user GET request
procedure GET(int c)
  cmsg m; m.c ← c;
  cmsg reply;
  res ⟨ans1, ans2⟩ ← GETLOC2(c);
  if ans1.id = myID or ans2.id = myID then return FETCH(c); // fetch from local disk if possible
  reply ← DOCALL(ans1.loc, REMOTEGET, m);
  if reply.resultCode = OK then return reply.content;
  // first successor failed; try the second successor
  reply ← DOCALL(ans2.loc, REMOTEGET, m);
  if reply.replycode = OK then return reply.content;
  // both successors failed; throw an error

```

**Figure 4:** Implementation of the content layer

**11. [8 points]:** What are the end-to-end properties of the content layer, assuming the algorithm in Figure 4? Assume that there are no failures of computers or disks while the system is running and that all tables are initialized correctly.

**(Circle True or False for each choice.)**

*There were two possible ways to answer this question, depending on whether the network was assumed to be faulty or fault-free. In the fault-free case, the answers are TRUE, TRUE, FALSE (for the same reason as (c) below), and TRUE (GET correctly finds any content stored with PUT). Here are the answers for the faulty case:*

- A. True / False** GET(*c*) always returns the same content that was stored with ID *c*.  
FALSE. *In the faulty case, the network could corrupt a packet in transit, in which case the content that was stored would be different than the content returned.*
- B. True / False** PUT(*c, content*) stores the content at the two immediate successors of *c*.  
FALSE. *The two immediate successors of c could be unavailable in the event of a network failure.*
- C. True / False** GET returns the content from the immediate successor of *c*.  
FALSE. *If the local node is not the immediate successor of c but is the second immediate successor of c, GET will return the id of the local node, as given by the fourth line in the pseudo-code for GET.*
- D. True / False** If the content has been stored on some computer, GET will find it.  
FALSE. *If the network has failed, some content may be unavailable.*

**12. [12 points]:** Now, suppose that there are failures while the system is running. Which of the following properties of the content layer are true?

**(Circle True or False for each choice.)**

- A. True / False** One of the computers returned by GETLOC2 might not answer the GET or PUT call.  
TRUE. *If one of the computers has crashed, it might not answer a call.*
- B. True / False** PUT will sometimes be unable to store the content at the content's two immediate successors.  
TRUE. *Both of the successors may have crashed.*
- C. True / False** GET will successfully return the requested content, assuming it was stored previously.  
FALSE. *The two nodes that previously stored the content may both have crashed.*
- D. True / False** If one of the two computers on which PUT the content is not failed when GET runs, GET will succeed in retrieving the content.  
EITHER TRUE OR FALSE. *If you assume that the computer that is not failed has not been partitioned from the network, this is TRUE. Partitioning would make this answer FALSE.*

### II.3 Improving Forwarding Performance

In this question we return to the forwarding and routing layer; you can ignore the content layer.

Alyssa isn't very happy with the performance of the system, in particular GETLOC. Her friend Lem E. Tweakit suggests the following change: each computer maintains a *NodeCache*, which contains information about the IDs and IP addresses of computers in the system. The *NodeCache* table initially contains information about the computers in *Successors*. The modified code for GETLOC is in Figure 5. You need not examine this code closely. It differs from the code in Figure 3 in only two ways: (1) it looks for the local answer in the *NodeCache* rather than in *Successors*, and (2) it inserts the new pair into the *NodeCache*. The changes are marked on the figure.

For example, initially the *NodeCache* at computer N232 contains entries for computers N251, N8, N21, N36, given the setup in Figure 1. But after computer N232 communicates with computer N36 and learns the ID and IP address of computer N96, its *NodeCache* now contains entries for computers N251, N8, N21, N36, and N96.

This code uses the same LOOKUP procedure shown in Figure 3. INSERT adds the new  $\langle id, loc \rangle$  pair to the *NodeCache* if it isn't already present, and it maintains the entries in the *NodeCache* in ring order; you may assume it does this correctly.

```

structure pair { int id; IP loc; }
structure res { int id; IP loc; bool succ };
structure msg { int c; int id; IP loc; bool found; int resultcode}
define succsSz 4;
pair Successors[succsSz];
pair NodeCache[1000]; // initially contains all entries in Successors
int NodeCacheSize; // number of entries in NodeCache; initially equals succsSz

// the procedure called by the user
procedure GETLOC(int c)
  res ans ← LOOKUP(NodeCache, NodeCacheSize, c); // check local NodeCache; for LOOKUP see Figure 3
  if ans.found then return ⟨ans.id, ans.loc⟩; // answer found locally
  msg m, reply;
  m.c ← c;
  while TRUE do { // search iteratively
    reply ← DOCALL(ans.loc, GETLOCLocal, m);
    if reply.resultCode = OK then {
      INSERT(NodeCache, ⟨reply.id, reply.loc⟩); // add response to NodeCache
      if reply.found then return ⟨reply.id, reply.loc⟩; // remote computer found answer
      ans.loc ← reply.loc; // remote computer returned its largest successor; continue searching there
    } else // error handling (not shown)
  }
}

// the handler invoked by a RPC request
procedure GETLOCLocal(msg req)
  msg reply;
  res ans ← LOOKUP(NodeCache, NodeCacheSize, c); // check local NodeCache; for LOOKUP see Figure 3
  reply.resultCode ← OK; reply.loc ← ans.loc; reply.id ← ans.id; reply.found ← ans.found;
  return reply;

```

**Figure 5:** Lem's implementation of GETLOC with caching

**13. [6 points]:** Assume that  $c$  is a content ID whose immediate successor is not one of the computers listed in *Successors*, and  $n$  is the number of computers in the system. In the *best case*, how many remote lookups are needed before  $\text{GETLOC}(c)$  returns?

(Circle the BEST answer)

- A. 0 *The correct answer is 0. In the best case, the NodeCache at the computer where GETLOC is called will contain the immediate successor of the requested content ID  $c$ . In this case, GETLOC will return the IP address of the immediate successor of  $c$  without having to make any RPCs.*
- B. 1
- C. 2
- D.  $O(\log n)$
- E.  $O(n)$
- F.  $O(n^2)$

**14. [12 points]:** Alyssa is wondering if the implementation suggested by Lem is correct. Which of the following statements are true about the implementation of figure 5, assuming no failures.

(Circle True or False for each choice.)

*This question was graded +3 for a correct answer, with no penalty for incorrect answers.*

- A. **True / False** If the immediate successor of  $c$  is present in *NodeCache*,  $\text{GETLOCLocal}$  will return this computer's address (*loc*).  
TRUE. *GETLOCLocal checks its NodeCache and if the immediate successor is there, it will return the IP address of that computer.*
- B. **True / False** If the immediate successor of  $c$  is not present in *NodeCache*, the reply message from  $\text{GETLOCLocal}$  will contain **found = false**.  
FALSE. *Unfortunately, Lem's implementation has a mistake in it. When the lookup procedure searches the NodeCache, it checks each entry, in ring order, starting from the machine on which it is running, to see if any of them has an ID that is greater than or equal to that of  $c$ . If it finds such a computer it returns a reply containing that computer's IP address and **found = TRUE**. The problem is that the NodeCache may contain incomplete information, and therefore the computer whose ID is returned in this response may not be the immediate successor of  $c$ . For example, the NodeCache at computer N232 might contain entries for computers N251, N8, N21, N36, and N96. Based on this information, if  $\text{getLocLocal}$  running at computer N232 is asked for the immediate successor of  $c = C50$ , it will return the IP address of N96, with **found = TRUE**, even though N96 isn't the immediate successor of C50.*
- C. **True / False** If the immediate successor of  $c$  is present in *NodeCache*,  $\text{GETLOC}$  will return this computer's address.  
TRUE. *This is true for the same reason that (A) is true.*
- D. **True / False**  $\text{GETLOC}$  will always return the immediate successor of  $c$ .  
FALSE. *This is false for the same reason that (B) is false.*

## End of Quiz II

Name: