582615: Overlay and peer-to-peer Networks, Autumn 2011

Exercise 3
Example solutions
(Questions which require pseudocode not included)

Questions

1. Answer the questions a) and b) based on the article "The Impact of DHT Routing Geometry on Resilience and Proximity "
http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.6.5914&rep=rep1&type=pdf

a) What are the main geometry types addressed in the article? Be prepared to draw and explain an example of routing in each geometry type on the blackboard at the exercise session. (3p)

Main geometry types addressed in the article are Tree , Hypercube , Butterfly , Ring , and XOR.

b) What are PNS and PRS? Why only certain geometries can utilize both PNS and PRS? (3p)

PNS means Proximity Neighbor Selection. In PNS the nodes to be inserted into  the routing table are chosen based on their proximity (measured in latency) when building a routing table.

PRS means Proximity Route Selection. PRS is applied when performing actual routing with the help of an existing routing table. In PRS the next-hop node is chosen from the already existing routing table based on the proximity (measured in latency) of the node.

Utilizing both PNS and PRS is not possible with some geometries, because some geometries don't allow proximity to affect Neighbor selection (e.g. butterfly) and others don't allow proximity to affect route selection (e.g. tree).

2. BitTorrent mainline DHT is a widely deployed Kademlia-based DHT with millions of concurrent users.

a) Explain how would you measure/estimate the number of concurrent users in BitTorrent mainline DHT. Note: take into account the vast size of the network. (3p)

There are $2^{160}$ possible node identifiers in the BitTorrent mainline DHT, which rules out the possibility of using brute-force approach of going through all node identifiers and checking if they are in use. Instead one needs to estimate the node density (or average distance between nodes) in the identifier space and estimate the total number of nodes based on this value.  Node density can be estimated by measuring the  node density in a number of smaller portions of the in identifier space.

The node IDs are random numbers generated by pseudo-random number generators, and there is churn in the DHT. This means that the node density is probably not evenly distributed in the identifier space. Thus it is necessary to sample the node density in several areas of the identifier space before jumping into conclusions.


There are several ways for estimating the node density and calculating the number of nodes in the DHT based on this information. Here one possible way:

First we estimate the average length of a node ID prefix which is needed to distinguish a node from another (average distance between nodes). We mark the length as l_unique. Then the number of concurrent nodes in the DHT is given by expression 2 ^ l_unique.

The prefix length can be approximated as follows. Find a set of consequtive nodes in terms of the XOR metric near a reference point. Calculate the mean of common prefix lengths of successive nodes in the set. This length is marked as l_common. Now an estimate for the unique prefix length is l_unique = l_common + 1.

This process needs to be repeated with several reference points and an average of the results must be taken, in order to overcome the effect of random variation.