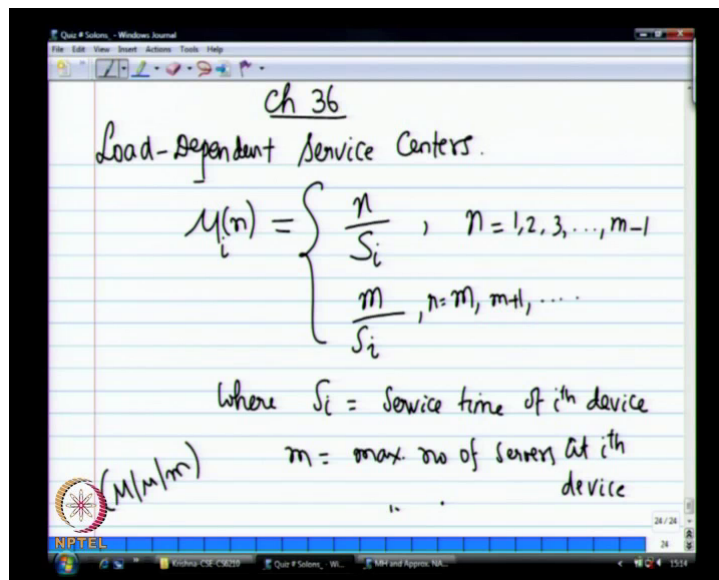


Performance evaluation of computer systems
Prof. Krishna Moorthy Sivalingam
Department of Computer Science and Engineering
Indian Institute of Technology, Madras
Lecture No. # 30
Load-Dependent Service Centres

(Refer Slide Time: 00:10)



So far we have only seen $m/m/1$ systems. So, how do we extend that analysis to systems where there are multiple servers in a given queue, serving a given queue.

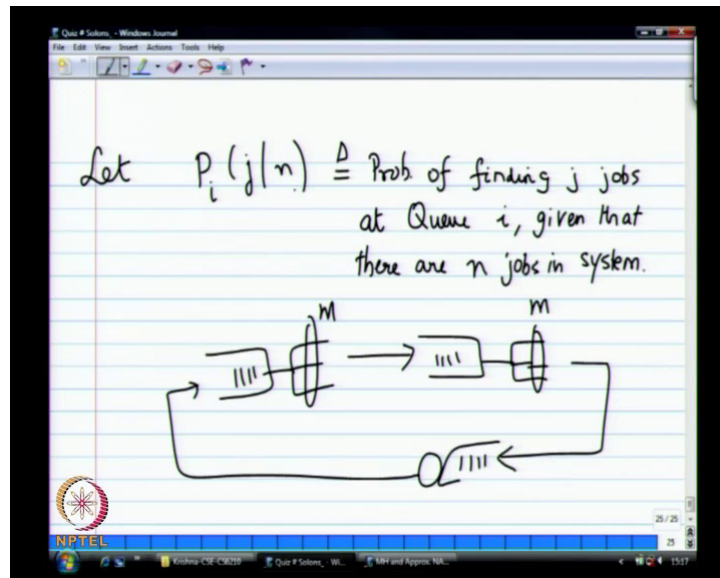
(No audio available: 00:27-00:44)

So, in a load dependent service center, previously we had the μ or S_i was simply just one over μ . Now, the service rate is going to depend upon the number of customers in the queue. So, if there is only one customer, the service rate is simply μ . We saw the $m/m/1$ system or $m/m/m$ system. If there are two customers, service rate is $2\mu, 3\mu$ and so on. Therefore, the effective capacity of the system changes with the number of customers in the queue. That is the only change which is not new to us. We know that.

So, the effective capacity of the system when there are n customers in the queue is defined as follows. So, if s is the service time, so it is simply one over s , right. The capacity is one over S

i. This case just as. You could let say, this is $(M/M/m)$. So, when the ith capacities when there are n customers in its queue is simply n by S_i for n upto m minus 1 m, where m is the number of servers, where S_i is the service time of the ith device and m is the maximum number of servers at ith device. Basically, your m/m/m server. So, this is other way of defining about the service of the queue varies with the number of customers waiting for service.

(Refer Slide Time: 03:05)



(No audio available: 2:48-3:10)

Now, couple of definitions $P_i(j/n)$. So, this is defined as the probability of finding j jobs at queue. So, if it m m 3 or 3 service in the ith device, so every device is now so what we are looking ok I will just come back. Now, I am defining a queue of networks where there are multiple servers available. Previously, only one server was there and then, collectively this feeds into another queue which could also be a multi server device and then, this feeds into another queue and so on which might be a single server device. So, you might have a system like this. So, there are m such servers. This is a collection of m/m/m queues.

(No audio available: 4:16-4:37)

So, this is $P_i(j/n)$ is probability of finding j jobs at queue i given that there are n jobs in the system. Again this is for a closed queueing network. So, with that definition, then we will try to define something.

(Refer Slide Time: 05:17)

The image shows a digital whiteboard with three mathematical formulas written in black ink on a lined background. The formulas are:

$$R_i(n) = \sum_{j=1}^n P_i(j-1 | n-1) \frac{j}{\mu_i(j)}$$

$$P_i(j | n) = \begin{cases} \frac{X(n)}{\mu_i(j)} P_i(j-1 | n-1), & j=1, 2, \dots, n \\ 1 - \sum_{k=1}^n P_i(k | n), & j=0 \end{cases}$$

$$Q_i(n) = \sum_{j=1}^n j P_i(j | n)$$

The whiteboard also features a toolbar at the top with various drawing tools and an NPTEL logo in the bottom left corner.

So, R_i which is the response time at server i , when there are n jobs in the system, it is given as follows.

(No audio available: 5:30-6:09)

Does it make sense? So, if there are j jobs at a given point in a time that is simply a probability of this j jobs, finding j minus 1 jobs when there are n minus 1 jobs in the system. This is $P_i(j-1 | n-1)$. n minus 1 is the probability of having j minus 1 jobs in the i th queue when the n minus 1 job is in the system. Then, your j th type job comes along. Therefore, the total service time is simply j divided by $\mu_i(j)$.

Makes sense? Some of you are saying yes, others are like what $I_i(n)$. Yeah, this is the average response time of i th device depending upon the number of customers in the device and there is only one customer or let us say when there are only two customers, the response time is given in this manner. So, the number of customers in that particular queue can be from one to n . In the system, yes. Out of which they could all be sitting in this queue or they could be sitting in some other queues. So, when there are j customers in the queue, the delay is given by $P_i(j-1 | n-1)$ into j divided by $\mu_i(j)$ is the number of customers because each customer required by $\mu_i(j)$ service time. j service time.

So, that is not that hard if you spend other couple of days on it if you want to. Now, this magical probability j of n . So, how do we compute this? Thus, once we compute this, everything else falls into place. So, this is actually were not deriving this (ρ) . This is defined as ρ of n by (ρ) . (No audio available: 8:07-8:14).

So, this is defined recursively.

(No audio available: 8:18-8:50)

So, this is well, I have to go back and dig up some of the other queuing network textbooks. Where is the connection to ρ by μ ? X is the number of x trooped of the system, right. So, this is the average number of customers completing the service at this queue. I do not have the intuitive explanation for that.

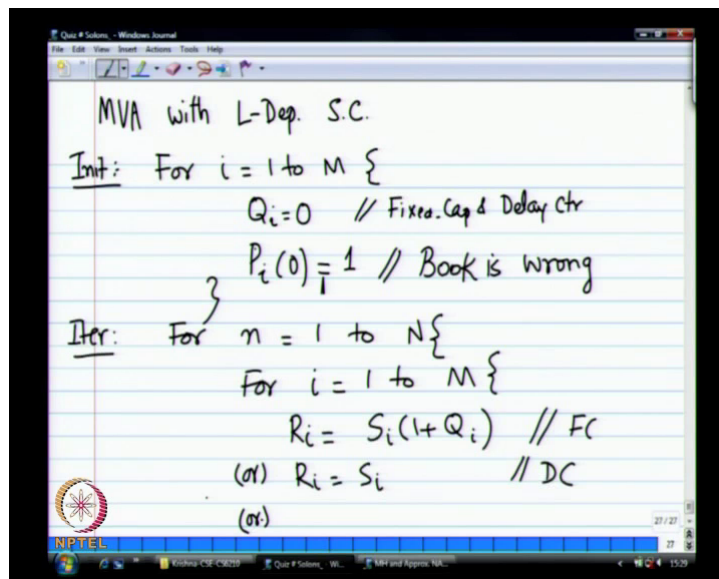
No, all it says is that if there is only one customer, let say the queue is there are 800mbps links, only one packet is in the queue, the capacity of the server is 100 mbps. When the two packets are simultaneously being served, the capacity is 400 or 200 mbps and so on. So, when all four servers are active, the capacity of the system is 400 mbps. No, it simply says that at the i th job, i th queue rather if there are n jobs, then the capacity of that server is n by S_i or m by S_i . If n is greater than or equal to m , then it becomes m by S_i . So, if I have 100 mbps links and there are four such links, right. If there is one packet in the queue, then the capacity of the system is how much? 100 mbps and that instant that is all it serving 100 m b p s. When there are two packets simultaneously being served, it goes to 200 mbps. When there are five packets in the queue, it is still 400 mbps because the others are getting queued. The effective capacity of the server is dependent upon the number of customers in the queue and so on.

So, you know it is 8, so it is going to be 400 mbps and then, Q_i of n . So, what will be Q_i of n ? You write them or can just tell me. What is the probability of i customers in this queue or when there are n customers in the system, what is the average number of customers at Q_i ? \sum_j . So, j goes from 1 to n and then, p_j . That is all. So, p_i tells me that what is the probability of the j customers of the n residing in Q_i and that into j will give me the overall Q_i . Now, we have expression for r and for q based on this definition of p .

Now, again we will have this (ρ) vary of enumerate for the different value of j as see go along. (ρ) . These are the expressions that are going to be changed with respect to the old MVA algorithm that we saw, exact MVA algorithm where we also compute. We compute at q ,

we compute at then r based on the q. Remember what we did? We said r equals s into one plus Q i and then, from the r, we found R i we found out the r, from that we found the x and from that Q is n did when repeated. So, in this MVA version, the same computation is being done. The only difference is that your computation of R is being done differently and computation of Q is being done differently. That is all. Otherwise, the same model that we saw, but the exact MVA approach. So, the algorithm itself is kind of you know similarly repetitive.

(Refer Slide Time: 13:00)



(No audio available: 12:49-12:58)

So, this is your MVA with so-called low dependent service content. Once I write the algorithm will so for (C). So, this is the exact MVA. We started by saying that Q i equal to 0 and then, we built up from Q i 1, Q i 2, Q i 3 and so on. So, this is for the fixed capacity and delay center. So, I have now a set of queues. Some are fixed capacity, some are delay center and some are this low dependent service center m/m/m. So, we love to run MVA for the entire systems and find out the trooped in other things. Now, we have to set this to one.

What is a probability of their being 0 customers? When there are 0 customers in the system, probability of 0 customers is Q i. When there are 0 customers in the system, it is simply one. That is the initialization. By the way the book is wrong. The book uses as 0 and it is incorrect. Start with 0; everything will be 0 if you look at the rest of the computation. So, that is an error that there the errata, but not in the chapter itself. So, make sure that you note this

down. Do not make this mistake when you come to the exams, you start with 0. You never will get 0 for everything else. This is your initialization. Then, your iteration.

(No audio available: 14:51-15:03)

So far when it is known for given a value, input value n, we compute Q. The input values for n equals 1 based on the n values equals 2 and so on. That is for exact MVA. So, this is now for every queue. So, what is R_i? If it is fixed capacity center, it was S_i into 1 plus Q_i. When there are n minus 1 customer, what is a queue length? Then, plus by self into the service time that so we calculated this. So, if it is for any i that is the fixed capacity center, this is R_i. So, look at this switch statement or k statement and R_i equals S_i if it is a delay center. So, depending on what i is representing, you would change Buff formula.

(Refer Slide Time: 16:18)

$$R_i = \sum_{j=1}^n P_i(j-1) \cdot \frac{j}{\mu_i(j)}$$

} // End for i=1 to M

$$R = \sum_{i=1}^M R_i V_i$$

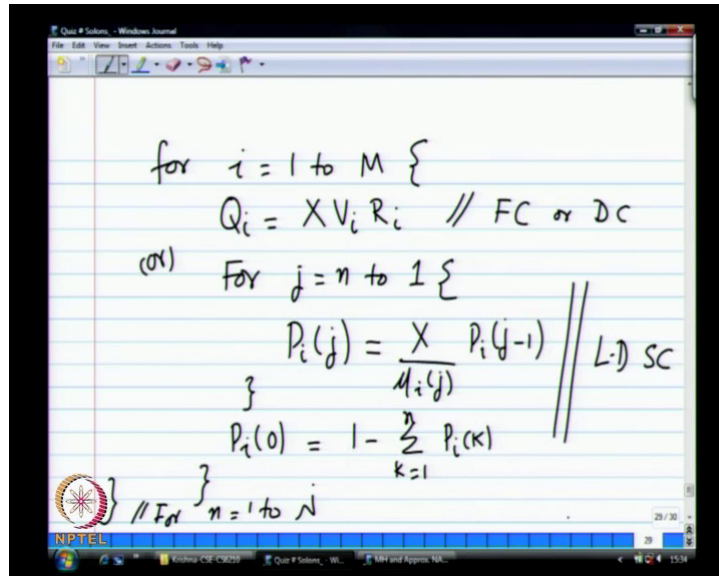
$$X = \frac{N}{R + Z}$$

(No audio available: 16:10-16:46)

So, n is your current value of the index, the outer book index. This is what the R_i equals when there are n customers side where they had the extra j minus 1 and then, if n minus 1 that if n minus 1 is what I removed. So, that is simply said I just need this 1, arrive p_i to do that. I do not need the two-dimension. Arrive it simply if there are j minus 1 customer in the system, then what that is into j. That is why they start with n equals 1, I need p_i of 0. P of 0 better be 1 that is 0. Then, I have everything else will be 0. So, that is why that book (()). So, this is it. So, this is my end for look. So, computed the R_i, then I compute my x. So, this is what

we did last time. Now, the next step is when we compute, there we recompute the Q, then will recompute the p i's based on the previous values.

(Refer Slide Time: 18:14)



So, remember how we computed Q i is equal to x into x i into R i. That is basically x into v i into R i. So, we recently calculated R i, x is also recently calculated. So, we update the values of Q i. This is basically Q i of 1. This is if it is a fixed capacity server or also delay center. Does not matter. So, here we have little bit of (0), so for j equals one to n. Sorry, there is one mistake in the book and which also has to be corrected. This is not because for n customers, this is again book will say book is in error. If you use that formula, that is you will be getting the wrong values of x where n is the original input, n is 20. That should not be used here. The current running value of the n variable is what should be used. Any other place where n appears. This hopes for i equals 1 to m. So, for a given value of n and computing all the R i's, computing the r, then the x and then, I am recomputing the Q i's to go to the next iteration of n.

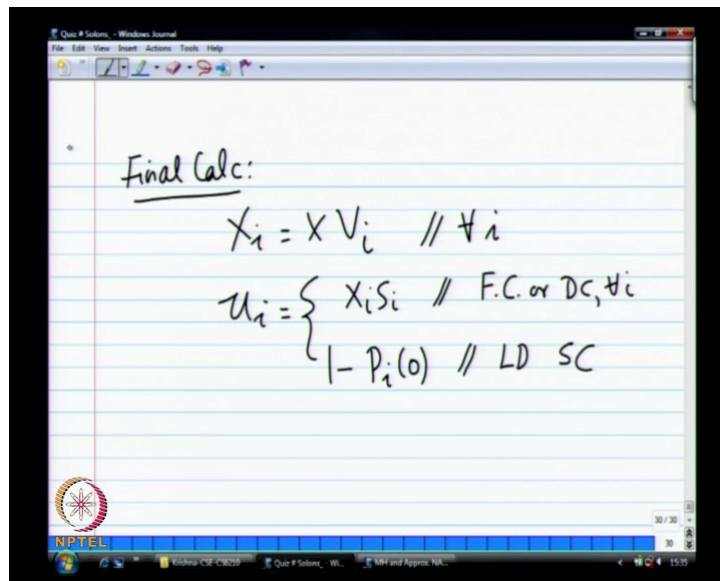
Smaller will go from 1 to n began is 20 will go from accurate MVA. We compute n equals one what are the q lines, n equals 2 what are the q lines and so on. So, this is for a given value of n within the, so for this we will again use that expression that we saw. This is the definition of p i of j, this is your x of n that we just computed divided by m u i of j into p i of j minus 1. Then, p i of 0 is simply 1 minus (No audio available:20:41-20:51).

So, this is for low dependent service.

(No audio available: 20:54-21:19)

So, there is no other place where upper case can appear good. So, that is one. This is one, this is one, this is as it is before. This is the only changed place on the p_i that we saw before and then, that loop is closed. Yeah, this closes the, yeah this is for n equals 1 to n . So, the end of n iteration you finally have the value (Q) . So, what we have is the x of 20 or r of 20 and so on.

(Refer Slide Time: 22:01)



The image shows a slide with handwritten mathematical equations. The title is "Final Calc:". Below it, the equations are:

$$X_i = X V_i \quad // \forall i$$
$$U_i = \begin{cases} X_i S_i & // \text{F.C. or DC, } \forall i \\ 1 - P_i(0) & // \text{LD SC} \end{cases}$$

The slide also features a logo in the bottom left corner and a status bar at the bottom with the text "NPTEL".

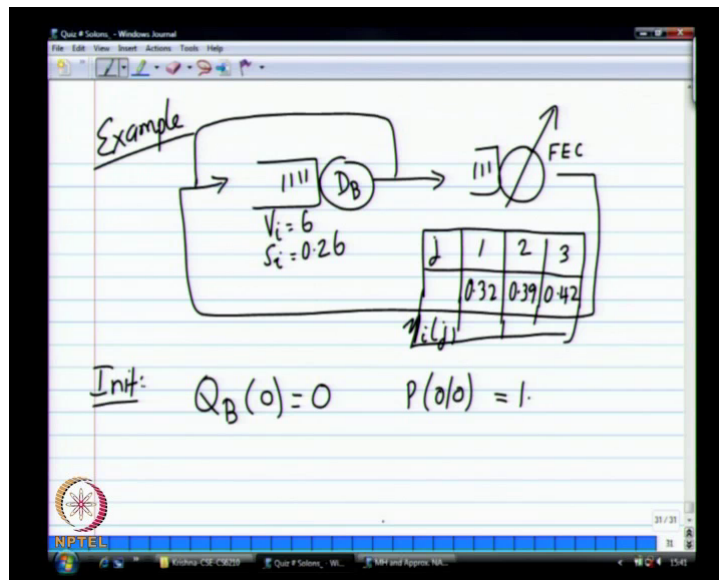
So, then we can do the final calculations. We will do x_i equals x into V_i . This is for all i and U_i is simply $X S$. So, $x_i S_i$ this is if it is a fixed capacity center or delay center for whatever values of i that match it and this is case of this one, it is simply (No audio available: 22:34-22:41).

So, p_i of 0 is the probability that there are no customers in the Q_i at the end of all the calculations and 1 minus that is simply going to be your utilization particular queue. Beginning of the equation are computing the R_i 's depending on whatever it is a service and fixed capacity or delay center or low dependent service center, right. That is this loopier. So, once I compute the R_i , I compute my r and x like you have done before. Then, once you compute my x , I go back and recompute the Q_i 's for all the devices. That is what this loop is all about.

Why it is a Q_i compare, yeah that this in matrix that is just for the sake of let see where we will use that. It is only to compute the final, yeah Q_i will be simply the number of customers that are in the thinking state, number of terminals that are in the thinking state. That is what a

Q i is. It is still Q i be simply U i. The Q i in that case will be equal to U i which is (∞) . Of course, we do not really have W i for that the delay center. Yeah, then delay center is where that R i W i is coming. Then, we have enough time to think for the numerical examples and will call it (∞) .

(Refer Slide Time: 24:31)



So, this is a shade of funny looking example. The reason for this will be to look at the next section. So, this is actually to serve the next section. So, we have this. This is your symbol for and we call this define this device FEC. So, jobs come from here to this disc and then, from the disc, it comes back here or it goes here. This is here. Only two queues. One is a fixed capacity queue. This is your fixed capacity center. This is your variable capacity center and definitions for this is here V_i is equal to 6 S_i equals 0.26. We saw the same fill up before, this disc b 0.6 and 0.26 and 6. What is happening here is the disc A and the CPU have been so-called collapsed, this variable capacity center. How we will do that? We will see later. Now, the definition for this queue, this is follows.

So, this is j and this is μ_j of j . So, when there is one customer, the service rate is 0.32. When it is 2, it is 0.39 and when it is 3, it is 0.4. It is little bit different from what we saw so far where it is always a multiple of the number of customers. Here it is not a multiple, but this is true for any value of μ . You can define what ever combinations. When expected value of number of customers increases, your capacity also increases. That is all. That is only real relationship.

In an m/m/1 system, m/m/m system as m, it is proportional to m and also note that everything defined for m greater than 4 which is also fine. So, this case cannot have more than four customers for whatever reason. Utmost three customers can be $(())$. So, this is our input. So, only two queues to worry about. This is a variable capacity server depending upon the number of customers in that queue.

So, the initialization parameters. So, remember Q_B , right. Q_B of 0 equals 0. When there are no customers, there are no, in the system itself there are no customers here and p of 0,0 equals 1. So, for a low dependent system, we use the probabilities p_j^n . This is 0 customers in this queue. When there are 0 customers in the system that is always equal to 1. This is our starting point.

(Refer Slide Time: 27:57)

Iter 1 : $n = 1$

$$R_B(1) = S_B (1 + Q_B(0))$$

$$= 0.26 (1 + 0) = 0.26$$

$$R_{FEC}(1) = P(0|0) \cdot \frac{1}{0.32} = 3.13$$

$$R = \sum R_i V_i =$$

So, we now go to the iteration one. So, n will take the value one. So, now there is only one customer in the system that is circulate. So, R_B of 1 is simply S_B into 1 plus Q_B . In the algorithm we did not have this parenthesis of 0 and so on, but this is what it is. So, this is 0.26 into 1 plus 0. Therefore, this is 0.26 and R_{FEC} of 1. Remember this is summation going from one to n . N this case happens to be one. So, this is simply P_{FEC} of p_0 . This FEC also will drop because this is the only device which has this problem. So, it is simply P_0 into j . So, j is 1 divided by 0.32. So, that is the first iteration. So, then you recompute R .

(No audio available: 29:22-29:37)

Oh yeah and neither input is that your, I said V. This is VV. Here your VFEC equals 1. Only basically.

(Refer Slide Time: 30:49)

$$Q_B(1) = X(1) R_B V_B$$

$$= 0.21 \times 26 \times 6 = 0.33$$

$$P(1|1) = \frac{X(1)}{M(1)} \cdot P(0|0) = \frac{0.21}{0.32} \times 1 = 0.67$$

$$P(0|1) = 1 - P(1|1) = 0.33$$

Iter 2: n=2 $R_B(2) = S_B (1 + Q_B(1)) = 0.26 (1 + 0.33)$

$$= 0.35$$

So, Q_B of 1, the updated version for the Q length is now based on the new x that we calculated. So, X into R_B into V_B . So, that is now 0.33. So, from Q_B of 0, it is now updated 0.33. Likewise p of 00 was 1. P of 11 is what based on p 00. So, that is now updated to 0.67 and if you look at the expression, we compute everything from one to n and p of 0 or p of that 0, that number of customers being 0 in the i th queue. It is simply 1 minus the probability of all the other problems, sum of all the other problems. So, this is the probability of one customer in that variable capacity center when there is only one customer in the system is 0.67. So, probability of no customer in this queue is when there is only one customer. The entire system is 0.3. This is time to compute. This is the first step of the iteration.

Then, we go to the second step. Now, my n gets the value 2. Then, we repeat this whole thing again. So, Q_B equals S_B into (Refer Slide Time: 30:56-30:59)

So, now Q_B with two customers. This is with Q_B with one customer. So, that is 0.26 into 1 plus 0.33. Oh sorry. Yeah, so this is R_B , but this is Q_B . Yes, now there Q_R and S_R are at proper place. Now, this is 0.35. So, this was the standard update for the last time. So, only now the probabilities are going to be different. Now, the probability of **(0)**.

(Refer Slide Time: 31:47)

$$R_{FEC}(2) = P(D|1) \cdot \frac{1}{\mu(1)} + P(1|1) \cdot \frac{2}{\mu(2)}$$

$$= 0.33 \times \frac{1}{0.32} + 0.67 \times \frac{2}{0.39} = 4.46 \text{ s}$$

$$R(2) = \sum R_i V_i = 6.54$$

$$X(2) = 2 / R(2) = 0.31$$

[j=1,2; n=2]

So, this is now RFEC with two customers equal to the probability of one customer or no customers in this queue with overall one customer into j by, so this is j by mu of j. Then, j takes the value. So, j is 1. This is j minus 1; this is n minus 1, j minus 1. So, in this case j is 2. Sorry, j goes from 1. j is either 1 or 2 and n equals 2. So, this is 0.33, this is 0.32, this was 0.67 into 2 divided by this was 0.39. So, now your mu r value is to be 4.46 seconds. So, it is kind of repetitive. That is the basic idea, but this helps you find what you are looking for. If you have (0). Now, RF2 is now computed like before R i V i and that is now updated to 6.54, then x of 2 is now 2 divided by R of 2 that is 0.31 and so on. Ok, that followed.

(Refer Slide Time: 34:09)

$$Q_B(2) = X(2) V_B R_B(2)$$

$$= 0.64$$

$$P(2|2) = \frac{X(2) \cdot P(1|1)}{\mu(2)}$$

$$= 0.52$$

$$P(1|2) = \frac{X(2) \cdot P(0|1)}{\mu(1)} = 0.32$$

$$P(0|2) = 1 - P(2|2) - P(1|2) = 0.16$$

Then, we do all the updates. So, Q_B of 2 is x into V_B into R_B of 2 or x of 2 and that is now x updated to 0.64. So, the queue length increases to 0.64. Now, we have to compute the P_{21} and then, P_{02} . So, P_{21} of 2 is given by x of 2 by μ of 2 P_{01} . So, this is 0.52. Then, P_{12} again. So, P_{12} of j is given by x of n which is that by μ of 1. **So, no more that.** Then, it is j minus 1 n minus 1. So, this is P_{01} and that is nothing, but 0.32. Now, what you find is the probability of 2. Both the customers being in this FEC when there are only two customers in the system, it is fairly high or 0.52 and then, one customer is 0.32 and therefore, no customers is simply $1 - P_{12} - P_{21}$. So, that should be 0.16. So, this is end of our second iteration.

(Refer Slide Time: 36:09)

The image shows a whiteboard with the following handwritten text:

$$\text{Iter 3 } R_B(3) = S_i(1 + Q_B(2)) = 0.2$$

$$R_{FEC}(3) = P(0/2) \cdot \frac{1}{\mu(1)} +$$

$$P(1/2) \cdot \frac{2}{\mu(2)} +$$

$$P(2/2) \cdot \frac{3}{\mu(3)}$$

$$= 5.86$$

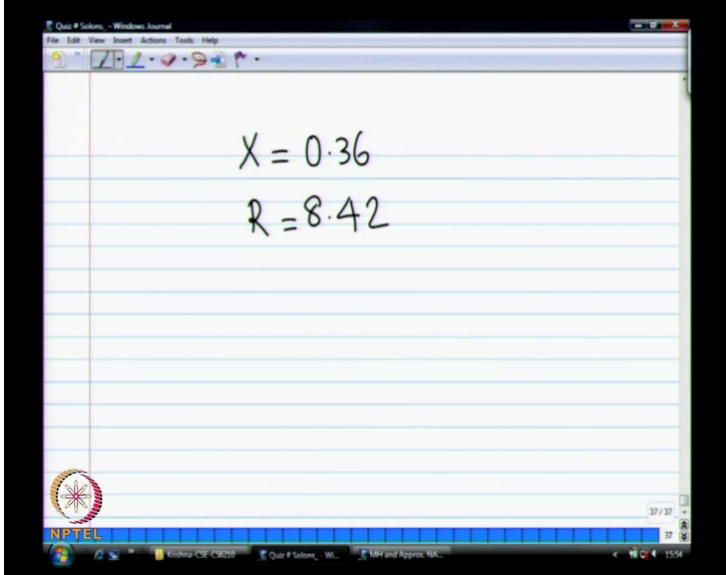
So, third iteration. I will let you figure it out. Now, R_{FEC} of 3 equals 1. So, again j goes from 1 to 3. Therefore, it is j divided by μ of j and then, $P_{j-1, n-j}$. So, if when j is 1, so this is 0.2 into 1 by μ of 1. When j goes to 2, then this is 2 by μ of 2. When j goes to 3, this is (Refer Slide Time: 36:47-36:51).

Yeah, this then finally computes to 5.86. Now, Q_B we know that is simple. Sorry, R_B , this is and so on. (Refer Slide Time: 37:08-37:26).

Once you get the r , you get the r . To get the R_i is get the r and then your X and then, you repeat this. We have one more iteration, but that we do not have to worry about, but in the end

you will be able to get the throughput of the system and so on. So, the final result we will just write that down.

(Refer Slide Time: 37:54)



A screenshot of a Windows Internet Explorer browser window. The window title is "Class # Solutions - Windows Internet". The address bar is empty. The main content area shows two handwritten equations: $X = 0.36$ and $R = 8.42$. The background of the page is lined paper. In the bottom left corner, there is a circular logo with a starburst pattern and the text "NPTEL" below it. The Windows taskbar is visible at the bottom, showing the Start button, several open applications, and the system tray with the time 11:54.

Now, your X is finally going to be 0.36 and your final R is going to be 8.42. So, we will stop here. Questions before we $(())$. So, this is the $(())$ that you can bring in there variable capacity servers also into the system. So, programming assignment three requires due to implement a bunch of $m/m/m$ server and solve that simulation and then, solve this very simple program and then, compare the two results.