This paper was appeared in IEEE Transactions on Computers Nol. C-29, no.7 (July 1980) pp. 611-617

# OPTIMIZATION OF MEMORY HIERARCHIES IN MULTIPROGRAMMED COMPUTER SYSTEMS WITH FIXED COST CONSTRAINT

Ъу

### TR 18-7

Samuel T. Chanson and Prem S. Sinha Department of Computer Science University of British Columbia

### ABSTRACT

This paper presents, using queuing theory and optimization techniques, a methodology for estimating the optimal capacities and speeds of the memory levels in a computer system memory hierarchy operating in the multiprogrammed environment. Optimality is with respect to mean system response time under a fixed cost constraint. It is assumed that the number of levels in the hierarchy as well as the capacity of the lowest level are known. The effect of the storage management strategy is characterized by the hit ratio function which, together with the device technology cost functions are assumed to be representable by power functions. It is shown that as the arrival rate of processes and/or the number of active processes in the system increase, the optimal solution deviates considerably from that under a uniprogrammed environment.

.

### Set a sec

### I. INTRODUCTION

The optimization of memory hierarchies is recognized as an important research area and has been attacked from several directions. Various solutions optimal under certain constraints have been obtained [1]-[5]. Ramamoorthy and Chandy [1] have obtained the size and type of each memory level by minimizing the average access time of an information block in a program for a given cost constraint. The concept is then extended to a general case in multiprogramming. The approach presupposes the knowledge of the frequency of access for each information block. MacDonald and Sigworth [3] have dealt with various combinations of optimization criteria such as fixed cost constraint, fixed and variable page size etc.. They too assume knowledge of the storage address sequence and have used its statistical properties extensively in their work. The objective function to be minimized is average access time or a function of it. Chow [2] has very nicely applied geometric programming to obtain not only the optimal size and speed of each memory level, but also the optimal number of levels for a given cost constraint. There it is assumed that storage management strategy is characterized by a hit-ratio function. Furthermore, the hit-ratio function and device technology cost function are taken as power functions of the capacity and access time respectively of each level of memory. Chow's analysis is restricted to uniprogrammed systems. Welch [4] gives a very simple and straightforward analysis of memory hierarchy for speed-cost trade-off with the assumption that the size and access probability of each level of memory are known and fixed. Rege [8] uses a simple two-server queuing network model to analyze the cost-performance trade-off by using different sizes and speeds at different memory levels. There is no optimization study.

All previous work in optimization use the mean memory access time as objective function for minimization and with the exception of [1], deal only with the uniprogrammed environment where only one process is active at any time and the processor simply hangs up when a request is made to any memory level. It is not clear that in a multiprogramming environment, where a process may be blocked while it is referencing information in certain memory level that minimizing the average memory access time is meaningful. In this paper we have combined performance evaluation techniques with optimization methods to extend the analysis of Chow [2] to cover multiprogrammed systems. Mean response time<sup>1</sup> is chosen as the objective function. With the number of memory levels fixed and the capacity of the lowest level known, we obtain from queuing theory an expression for mean system response time in terms of the capacity and speed of each memory level. The optimal expression of memory sizes and speeds are then obtained using the Lagrangian function under a cost constraint.

Notice that in the uniprogramming environment,

Average response time =  $c_1 * Average$  access time of the

memory hierarchy + c<sub>2</sub>

where  $c_1 = average number of accesses to the memory$ 

hierarchy per interaction

and c2 = mean CPU time demand of the process per interaction

c<sub>1</sub> and c<sub>2</sub> are constants for a given process. Hence average response time and average memory hierarchy access time are equivalent objective functions in the uniprogramming environment.

- 2 -

<sup>&</sup>lt;sup>1</sup> The term response time may be interpreted as request completion time in this paper.

#### II. SYSTEM DESCRIPTION, ASSUMPTIONS AND NOTATIONS

The memory hierarchy consists of N levels,  $M_1, M_2, \ldots, M_N$ , where N is known and fixed. Generally, the higher the level (i.e., the smaller the index) the smaller is the capacity, the faster its speed and the more expensive is its unit cost. It is assumed that information present in any level is also present in all subsequent lower levels. In the case of a uniprogrammed system, whenever the needed information is not found in the highest level M1, a request is made to each of the lower level successively until it is found in a level M,, i = 2, ..., N. The processor is held waiting all the time until the information is retrieved from M<sub>i</sub>. As i increases, the time required to fetch the information goes up. When i exceeds a certain value, it becomes uneconomical to keep the processor idle while the information is being retrieved from M,, particularly when there are other processes waiting for the processor. Thus in the case of multiprogramming, we have two types of memory - A and B. While the processor waits for access to type A memory, it does not do so for access to type B memory, but releases the current process and takes the next process ready to run if one exists. It is therefore possible for several requests to queue up at a type B memory level but there is at most one request at anyone time for a type A memory level. The model of such a system is shown in Figure 1 where  $n_1$  and  $n_2$  are the number of type A and type B memory levels respectively.  $X_i$ , i = 1, 2, ..., N (N =  $n_1 + n_2$  = a known constant) is the capacity of memory level M<sub>i</sub> in the hierarchy.  $X_N$  is assumed to be known. In addition, in consistent with Chow's [1] terminology, we define the following:

 $y_i$ , i = 1, 2, ..., N is the mean transfer time of a unit of information from level  $M_i$  to  $M_{i-1}$  (this does not include the queue wait time for type B memory levels) and  $y_1$  is simply the mean access time of the fastest memory.

- 3 -



FIGURE 1: STORAGE HIERARCHY IN A MULTIPROGRAMMED SYSTEM

H(x) is the probability of finding the requested information in a memory level with capacity x.

The hit-ratio function  $p_i$  (i.e., the probability of successfully retrieving the needed information from level  $M_i$ ) is therefore given by the difference in the probability of finding the information in  $M_i$  but not in  $M_{i-1}$ .

i.e., 
$$p_i = H(x_i) - H(x_{i-1})$$
,  
 $i = 1, 2, ..., N; H(x_0) = 0$  (1)

The missing ratio F(x) is simply 1 - H(x) and is assumed to be a power function of capacity x and positive constants  $K_1$  and  $\alpha$ , defined as

$$F(x) = K_1 x^{-\alpha}$$
 (2)

The technology cost function (i.e., unit cost of a storage level with transfer time y to the next higher level) is assumed to take the form

$$b(y) = K_2 y^{-\beta}$$
(3)

where  $\beta$  and K<sub>2</sub> are positive constants. Without loss of generality, we take K<sub>1</sub> = K<sub>2</sub> = 1, i.e., equations (2) and (3) become

 $F(x) = x^{-\alpha}$  (2')

and  $b(y) = y^{-\beta}$  (3')

This means that  $K_1^{\overline{\alpha}}$  is the unit for storage capacity and  $K_2$  is the unit for cost. Empirical data have shown that equations (2) and (3) are good approximations (see [2], [4]).

- 5 -

# III. QUEUING MODEL

The queuing model of the system described in the previous section is shown in Figure 2.



# FIGURE 2: QUEUING NETWORK MODEL A OF SYSTEM

where  $\omega$  is the mean rate of arrival of tasks (or requests in an interactive environment) to the system and is assumed to have an exponential distribution.  $q_1, q_2, \ldots, q_{n_1}$  are the probabilities of referencing memory levels  $M_1, M_2, \ldots, M_{n_1}$  respectively and  $p_1, p_2, \ldots, p_{n_2}$  are the probabilities of referencing memory levels  $M_{n_1+1}, \ldots, M_{n_1+n_2}$  respectively. The probability of exit (i.e., termination of task or completion of a request) is p.

Define

 $Q = \sum_{i=1}^{n_1} q_i ,$   $P = \sum_{i=1}^{n_2} p_i + p_i ,$ 

then clearly, P + Q = 1

For type A memory, the mean effective hierarchy access time  $T_i$  to level  $M_i (i \le N_1)$  is the sum of the mean individual transfer time between two consecutive levels from  $M_1$  up to  $M_i$ 

i.e., 
$$T_i = \sum_{j=1}^{i} y_j$$

In the case of type B memory the mean transfer times  $y_i$ 's are taken as the inverse of the service rates of the memory levels. Furthermore, it is assumed that the service rates are exponentially distributed with mean  $1/y_i$ ,  $i = n_1 + 1, \dots, n_1 + n_2$ . The mean effective hierarchy access time of type B memory levels are not so easily obtained because of possible queuing of requests at these levels. Furthermore, since a process may be blocked while accessing these levels it is more reasonable to use mean response time (or mean request completion time) as the criteria of optimization. To do so, we first transform the model in Figure 2 to the model in Figure 3. It is easy

- 7 -

to show, using first principles in queuing theory that the two models are equivalent with respect to the objective function.



# FIGURE 3: QUEUING NETWORK MODEL B OF SYSTEM

$$\mu' = P K$$

with 
$$1/K = \sum_{i=1}^{n_1} q_i/\mu_i$$

- 9

 $\mu_i = 1/T_i$  i = 1, 2, ...,  $n_1$  (4)

Taking  $\omega + \omega_1, \omega_1, \omega_2, \ldots, \omega_{n_2}$  to be the arrival rates of centres  $c_0, c_1, \ldots, c_n$  respectively, the mean response time of this network is given by [6]

$$\overline{R} = \frac{1}{\omega} \left[ \sum_{i=1}^{n_2} \frac{\omega_i \, y_{i+n_1}}{1 - \omega_i \, y_{i+n_1}} + \frac{(\omega + \omega_1)/\mu'}{1 - (\omega + \omega_1)/\mu'} \right]$$
(5)

$$= \frac{\sum_{i=1}^{n_{1}} x_{i-1}^{-\alpha} y_{i}}{1 - \sum_{i=1}^{n_{1}} \mu x_{i-1}^{-\alpha} y_{i}} + \sum_{i=n_{1}+1}^{n_{1}+n_{2}} \frac{x_{i}^{-\alpha} y_{i}}{1 - \mu x_{i}^{-\alpha} y_{i}}$$
(5')

where  $\mu = \omega/p$  (see Appendix 1 for details of derivation).

### IV. OPTIMIZATION

Since the technology cost per unit information is given by

$$b(y) = y^{-\beta}$$

the system cost with storage sizes  $x_1, x_2, \ldots, x_N$  for levels  $M_1, M_2, \ldots, M_N$ with average transfer times  $y_1, y_2, \ldots, y_N$  respectively, is given by

$$S = \sum_{i=1}^{N} x_i y_i^{-\beta}$$

(6)

Given N,  $x_N$ , and that the memory system cost is not to exceed  $S_0$ , the optimi- zation problem becomes:

$$Min \frac{\sum_{i=1}^{n_{1}} x_{i-1}^{-\alpha} y_{i}}{1 - \sum_{i=1}^{n_{1}} \mu x_{i-1}^{-\alpha} y_{i}} + \sum_{\substack{i=n_{1}+1 \\ i=n_{1}+1}}^{n_{1}+n_{2}} \frac{x_{i-1}^{-\alpha} y_{i}}{1 - \mu x_{i-1}^{-\alpha} y_{i}}$$
s.t. 
$$\sum_{\substack{i=1 \\ i=1}}^{N} x_{i} y_{i}^{-\beta} \leq S_{0} ,$$

$$x_{0} = 1 ; x_{i} > 0 ; y_{i} > 0 \qquad i = 1, 2, ..., N \qquad (7)$$

The problem (7) will have a solution only in the region where

$$\sum_{i=1}^{n_{1}} \mu x_{i-1}^{-\alpha} y_{i} < 1$$
  
and  $\mu x_{i-1}^{-\alpha} y_{i} < 1$   $i = n_{1}+1, \dots, n_{1}+n_{2}$ 

This restriction meets one of the assumptions made while calculating equilibrium state probability, that the traffic intensity has to be strictly less than one [6].

Now by multiplying the objective function by  $\mu$  (a constant) and adding  $1 + n_2$  to it, the problem (7) reduces to

$$\begin{array}{l} \text{Min} \quad \frac{1}{1 - \sum\limits_{i=1}^{n_{1}} \mu \, x_{i-1}^{-\alpha} \, y_{i}} + \sum\limits_{i=n_{1}}^{n_{1}+n_{2}} \frac{1}{1 - \mu \, x_{i-1}^{-\alpha} \, y_{i}} \\
\text{s.t.} \quad \frac{1}{s_{0}} \, \sum\limits_{i=1}^{N} \, x_{i} \, y_{i}^{-\beta} \leq 1 
\end{array} \tag{8}$$

The natural constraints  $x_i > 0$  and  $y_i > 0$  can be ignored in the calculation by looking at a solution only in the positive region of x and y.

Introducing new variables  $r_0$  and  $r_i$ 's such that

$$r_0 \leq 1 - \sum_{i=1}^{n_1} \mu x_{i-1}^{-\alpha} y_i$$

and  $r_{i-n_1} \leq 1 - \mu x_{i-1}^{-\alpha} y_i$ ;  $i = n_1 + 1, \dots, n_1 + n_2$ 

The problem (8) is equivalent to

We will use the standard Lagrange multiplier method together with geometric programming techniques to solve problem (9). The Lagrangian function for problem (9) is:

$$F(R,X,Y,\lambda) = r_{0}^{-1} + \sum_{i=1}^{n_{2}} r_{i}^{-1} + \lambda_{0} (r_{0} + \sum_{i=1}^{n_{1}} x_{i-1}^{-\alpha} y_{i} - 1) + \sum_{i=n_{1}+1}^{n_{1}+n_{2}} \lambda_{i-n_{1}} (r_{i-n_{1}} + \mu x_{i-1}^{-\alpha} y_{i} - 1) + \lambda' (\frac{1}{S_{0}} \sum_{i=1}^{N} x_{i} y_{i}^{-\beta} - 1)$$
(10)

Differentiating with respect to R, X, Y and  $\lambda$  respecting and equating to zero, we obtain,

$$\frac{\partial F}{\partial r_0} = (-1) r_0^{-1} + \lambda_0 r_0 = 0$$
(11)

$$\frac{\partial F}{\partial r_{i}} = (-1) r_{i}^{-1} + \lambda_{i} r_{i} = 0 ; i = 1, 2, ..., n_{2}$$
(12)

$$\frac{\partial F}{\partial x_{i}} = \begin{cases} (-\alpha) \ \mu \ \lambda_{0} \ x_{i}^{-\alpha} \ y_{i} + \frac{\lambda'}{S_{0}} \ x_{i-1} \ y_{i-1}^{-\beta} = 0 \quad ; \quad i = 1, 2, ..., n_{1} \end{cases}$$
(13)

$$\left((-\alpha) \ \mu \ \lambda_{i-n_{1}} \ x_{i-1}^{-\alpha} \ y_{i} + \frac{\lambda'}{s_{0}} \ x_{i-1} \ y_{i-1}^{-\beta} = 0 \quad ; \quad i = n_{1}+1, \ \dots, \ n_{1}+n_{2} \quad (14)$$

$$\frac{\partial F}{\partial y_{i}} = \begin{cases} \mu \lambda_{0} x_{i-1}^{-\alpha} y_{i} + (-\beta) \frac{\lambda'}{S_{0}} x_{i} y_{i}^{-\beta} = 0 ; i = 1, 2, ..., n_{1} \end{cases}$$
(15)

$$\left(\mu \lambda_{i-n_{1}} x_{i-1}^{-\alpha} y_{i} + (-\beta) \frac{\lambda'}{s_{0}} x_{i} y_{i}^{-\beta} = 0 ; i = n_{1}+1, \dots, n_{1}+n_{2} \right)$$
(16)

$$\frac{\partial F}{\partial \lambda_0} = r_0 + \sum_{i=1}^{n_1} \mu x_{i-1}^{-\alpha} y_i - 1 = 0$$
(17)

$$\frac{\partial F}{\partial \lambda_{i-n_{1}}} = r_{i-n_{1}} + \mu x_{i-1}^{-\alpha} y_{i} - 1 = 0 ; \quad i = n_{1}+1, \dots, n_{1}+n_{2}$$
(18)

$$\frac{\partial F}{\partial \lambda'} = \frac{1}{s_0} \sum_{i}^{N} x_i y_i^{-\beta} - 1 = 0$$
(19)

Define 
$$f^{0} = \lambda_{0} \sum_{i=1}^{n_{1}} \mu x_{i-1}^{-\alpha} y_{i} + \sum_{i=n_{1}+1}^{n_{1}+n_{2}} \lambda_{i-n_{1}} \mu x_{i-1}^{-\alpha} y_{i}$$
 (20)

$$\delta_0 = r_0^{-1} \tag{21}$$

$$\delta_i = r_i^{-1}$$
  $i = 1, 2, ..., n_2$  (22)

$$\delta_{1i} = \int_{0}^{\lambda_{0}} \mu x_{i-1}^{-\alpha} y_{i}^{\prime} f^{0} \qquad i = 1, 2, ..., n_{1}$$
(23)

$$\begin{cases} \lambda_{i-n_{1}} \mu x_{i-1}^{-\alpha} y_{i}^{\prime} f^{0} & i = n_{1} + 1, \dots, n_{1} + n_{2} \end{cases}$$
(24)

- 12 -

$$\delta_{20} = \lambda_0 r_0 \tag{25}$$

$$i = \lambda_i r_i$$
  $i = 1, 2, ..., n_2$  (26)

$$\delta_{3i} = \frac{\lambda'}{s_0 f^0} x_i y_i^{-\beta} \qquad i = 1, 2, ..., n_1 + n_2$$
(27)

Now clearly,

δ2

$$\sum_{i=1}^{n_1+n_2} \delta_{1i} = 1 \quad (normality) \tag{28}$$

$$(11) \Rightarrow (-1) \ \delta_0 + \delta_{20} = 0 \tag{29}$$

$$(12) \Rightarrow (-1) \delta_{1} + \delta_{21} = 0 ; i = 1, 2, ..., n_{2}$$
 (30)

(13) 
$$\&$$
 (14)  $\Rightarrow$  (- $\alpha$ )  $\delta_{1i} + \delta_{3i-1} = 0$ ;  $i = 2, ..., n_1 + n_2$  (31)

(15) 
$$\&$$
 (16)  $\Rightarrow$   $\delta_{1i} + (-\beta) \delta_{3i} = 0$ ;  $i = 1, 2, ..., n_1 + n_2$  (32)

Solving (28), (31) and (32) simultaneously, we obtain

$$\delta_{1i} = (\alpha\beta)^{N-i} \frac{\alpha\beta - 1}{(\alpha\beta)^N - 1}$$
  

$$\delta_{3i} = \frac{1}{\beta} \delta_{1i}$$
  

$$N = n_1 + n_2 \quad ; \quad i = 1, 2, ..., n_1 + n_2$$
(33)

Now in order to obtain optimal values for  $x_i$ 's and  $y_i$ 's we first obtain values for  $f^0$  and  $\lambda_i$ 's.

Raising (23) and (24) to the power  $\delta_{1i}$  and (27) to the power  $\delta_{3i}$  for  $i = 1, 2, ..., n_1 + n_2$  respectively, and then multiplying we obtain

$$(f^{0})^{1} \overset{\sum \delta_{1i}}{=} f^{0} = \lambda_{0}^{a_{n_{1}}} \underset{i=1}{\overset{n_{2}}{=} 2} \lambda_{i}^{b_{i}} c$$
 (34)

- 13 -

where 
$$a_{n_1} = \sum_{i=1}^{n_1} \delta_{1i}$$
  
 $b_i = \delta_{1n_1+i}$ ,  $i = 1, 2, \dots, n_2$   
 $c = \mu(\frac{1}{\beta S_0})^{1/\beta} \cdot x_N \cdot \prod_{i=1}^{n_1+n_2} (\frac{1}{\delta_{1i}})^{\delta_{1i}} (\frac{1}{\delta_{3i}})^{\delta_{3i}}$   
 $\frac{\lambda^i}{f^0} = \frac{1}{\beta}$ 

from (17), (21), (23), (25) and (29)

$$\lambda_0 = (2f^0 a_{n_1} + \sqrt{4f^0 a_{n_1} + 1})/2$$
(35)

and from (18), (22), (24), (26) and (30)

$$\lambda_{i} = (2f^{0} b_{i} + 1 + \sqrt{4f^{0} b_{i} + 1})/2$$
(36)

Substituting (35) and (36) into (34)

$$f^{0^{*}} = c \cdot (2f^{0} a_{n_{1}} + 1 + \sqrt{4f^{0} a_{n_{k}} + 1})^{a_{n_{1}}}$$
$$\cdot \prod_{i=1}^{n_{2}} (2f^{0} b_{i} + 1 + \sqrt{4f^{0} b_{i} + 1})^{b_{i}}/2$$
(37)

Substituting the value of  $f^{0^*}$  in (35) and (36) we obtain values for the  $\lambda_i$ 's.

From equations (23), (24), (27), (31), and (32), with some algebraic manipulation, we obtain the following expression for optimal  $x_i$ 's and  $y_i$ 's

- 14 -

$$\log x_{i} = \left[\frac{i}{1-a} - \frac{(1-a^{1})N}{(1-a)(1-a^{N})}\right] \cdot K + \frac{1-a^{1}}{(1-a^{N})} \cdot \log x_{N} + \frac{1-a^{1}}{(1-a)} \cdot \frac{L-M}{(1-a^{N})} \qquad i = 1, 2, ..., n_{1-1}$$
(38a)

$$\log x_{i} = \left[\frac{i}{1-a} - \frac{(1-a^{i})N}{(1-a)(1-a^{N})}\right] \cdot K + \frac{1-a^{i}}{(1-a^{N})} \cdot \log x_{N} + \frac{1-a^{i}}{1-a}\left[\frac{L-M}{1-a^{N}}\right] + \sum_{j=0}^{1-n} h_{j} a^{i-j} \qquad i = n_{1}, \dots, n_{1}+n_{2-1}$$
(38b)

where

 $a = \alpha \beta$ 

$$K = -(\beta + 1) \log a$$

$$h_{i} = \beta(\log \lambda_{i} - \log \lambda_{i+1})$$

$$L = \sum_{i=0}^{n_{2}-1} h_{i}$$

$$M = \sum_{i=0}^{n_{2}-1} h_{i} a^{n_{2}-i}$$

Also from (27) and (33), we have

 $y_{i} = (x_{i}/s_{0}, \delta_{1i})^{1/\beta}$  (39)

Hence knowing all  $x_i$ 's, all  $y_i$ 's can be found from equation (39). In the case of uniprogramming,  $\mu = 0$ , all  $\lambda_i$ 's are equal and therefore  $h_i = L = M = 0$ . Equations (38a) and (38b) then reduce to

$$\log x_{i} = \left[\frac{i}{1-a} - \frac{(1-a^{i})N}{(1-a)(1-a^{N})}\right] \cdot K + \frac{1-a^{i}}{1-a^{N}} \log x_{N} \qquad i = 1, 2, ..., N \quad (40)$$

This is the same result obtained by Chow [2]. The following example illustrates how the optimal solutions for a multiprogrammed environment approach the ones for a uniprogrammed environment as  $\mu$  tends to zero.

## Example

For the case N = 4,  $n_1 = n_2 = 2$ , and  $\alpha = \beta = 1$ , the expression for optimal  $x_i$ 's are

$$x_{1} = x_{4}^{1/4} \left(\frac{\lambda_{0}}{\lambda_{1}}\right)^{1/2}$$
$$x_{2} = x_{4}^{1/2} \left(\frac{\lambda_{0}}{\lambda_{1}}\right)^{2}$$
$$x_{3} = x_{4}^{3/4} \left(\frac{\lambda_{0}}{\lambda_{1}}\right)^{5/2}$$

Taking  $x_4 = 10^8$  and  $S_0 = 4 * 10^{10}$ , (actual units depend upon the normalization factors used in Equations (2) and (3)) the values of  $x_i$ 's are given in Table 1 for different values of  $\mu$ .

Table 1

Optimal storage size for N = 4,  $x_4 = 10^8$ ,  $S_0 = 4 \times 10^{10}$ ,  $\alpha = \beta = 1$ 

μ	x1	x <sub>2</sub>	×3
.100	100.08	10031	1003905
1000	100.77	10313	1039305
10000	107,24	13224	1418069
0	100	10000	1000000

The values of x's for  $\mu = 0$  are the ones obtained by Chow for a uniprogrammed system (the values are obtained using equation (40)). We observe that there is a significant difference when  $\mu(= w/p)$  is large (i.e., when the arrival rate of requests is large and/or the exit probability of a process is small which implies a large number of active processes competing for limited systems resources simultaneously). Also from Equations (38a), (38b) and (39), the ratio between the sizes and speeds of successive memory levels of the <u>same type</u> are constant, as observed by Chow and reported in [9].

#### CONCLUSION

We have presented a methodology for estimating the optimal capacities and speeds of the memory levels in a memory hierarchy operating in the multiprogrammed environment. Optimality is with respect to mean system response time under a fixed cost constraint. Mean response time is the single performance parameter most important to users of interactive systems which are invariably multiprogrammed. It is therefore more meaningful in such an environment to use mean response time rather than mean memory access time as the criterion for optimization. Queuing theory allows us to analyze multiprogrammed systems and global performance indices such as mean response time while optimization techniques enable us to compute the optimal values. This model can be further refined by relaxing certain assumptions made in this paper at the expense of mathematical tractability. For example by assuming a to be constant for all jobs (or requests), we are assuming that they have similar memory access characteristics. This assumption can be removed by taking  $\alpha_{ij}$  as the hit-ratio constant for the i<sup>th</sup> job in the j<sup>th</sup> memory level.

Since an explicit closed form solution for  $f^0$  does not seem to exist, we are unable to obtain optimal values for  $n_1$  and  $n_2$  as has been done by Chow [2] for uniprogrammed systems.

# ACKNOWLEDGEMENTS

This work was supported in part by the National Research Council of Canada under Grant No. A3554 and by UBC Summer Session Research Scholarship 4904.

#### REFERENCE

- C. V. Ramamoorthy and K. M. Chandy, "Optimization of memory hierarchies in multiprogramming system." J. Assoc. Comput. Mach., Vol. 17, pp. 426-445, July, 1970.
- [2] C. K. Chow, "On optimization of storage hierarchies." IBM J. Res. Develop., Vol. 18, pp. 194-203, May, 1974.
- [3] J. E. MacDonald and K. L. Sigworth, "Storage hierarchy optimization procedure." IBM J. Res. Develop., Vol. 19, pp. 133-140, March, 1975.
- [4] T. A. Welch, "Memory hierarchy configuration analysis." IEEE Trans. on computers, Vol. C-27, No. 3, pp. 408-413, May, 1978.
- [5] R. L. Maltson, "Evaluation of multilevel memories." IEEE Trans. Magnetics, Vol. MAG-7, pp. 814-819, December, 1971.
- [6] D. Ferrari, Computer system performance evaluation. Prentice-Hall, Englewood Cliffs, New Jersey, 1978.
- [7] R. J. Duffin, E. L. Peterson and C. Zener, <u>Geometric Programming</u>. John Wiley & Sons Inc., New York, 1967.
- [8] S. L. Rege, "Cost, performance and size trade-off for different levels in memory hierarchy." Computer, Vol. 19, pp. 43-51, April, 1976.
- [9] Y. S. Lin and R. L. Mattson, "Cost-performance evaluation of memory hierarchies." IEEE Trans. Magnetics MAG-8, pp. 390-392, Sept. 1972.

# Appendix 1

The arrival rates for service centres are

$$\begin{split} \omega_{n_{2}} &= (\omega + \omega_{1}) p_{n_{2}}^{*}, \\ \omega_{n_{2}-1} &= (\omega + \omega_{1}) p_{n_{2}-1}^{*} + \omega_{n_{2}} \\ &= (\omega + \omega_{1}) (p_{n_{2}-1}^{*} + p_{n_{2}}^{*}), \\ \omega_{1} &= (\omega + \omega_{1}) \sum_{j=1}^{n_{2}} p_{j}^{*}, \\ &\vdots \\ \omega_{1} &= (\omega + \omega_{1}) \sum_{j=1}^{n_{2}} p_{j}^{*}, \end{split}$$

and

Let

then

Now  $\omega + \omega_1 = \omega + \frac{P^*}{Q^*} \omega = \frac{\omega}{Q^*}$ 

⇒

(A1.2)

(A1.1)

. from (A1.1)  $\omega_i = \frac{\omega}{Q'} \begin{bmatrix} n_2 \\ \sum p'_j \end{bmatrix}$ ,  $\omega_{i} = \frac{\omega}{p} \sum_{j=i}^{n_{2}} p_{j}$ or

 $P' = \sum_{j=1}^{n_2} p'_j$ 

Q' = 1 - P',

 $\omega_1 = \frac{P'}{Q'} \omega$ 

(A1.3)

because Q' = 1 - P'

$$= 1 - \sum_{1}^{n_{2}} p'_{j}$$

$$= 1 - \frac{1}{P} (\sum_{1}^{n_{2}} p_{j})$$

$$= 1 - \frac{1}{P} [P - p]$$

$$= \frac{P}{P}$$

From Equation (1)

$$p_{j} = H(x_{U}) - H(x_{U-1}) \qquad U = j + n_{1}$$

$$\sum_{j=1}^{n} p_{j} = \sum_{j=1}^{n_{1}+n_{2}} [H(x_{j}) - H(x_{j-1})]$$

(taking 
$$H(x_N) = 1$$
)  
= 1 -  $H(x_{i-1})$   
=  $F(x_{i-1})$   
=  $x_{i-1}^{-\alpha}$ 

substituting in (A1.3) we get

$$\omega_{i} = \frac{\omega}{p} x_{i-1}^{-\alpha}$$

$$= \mu x_{i-1}^{-\alpha}$$
(A1.

9

where

 $\mu = \frac{\omega}{p}$ 

.4)

Also

$$(+ \omega)/\mu' = \frac{\omega}{Q'\mu'}$$
$$= \frac{\omega}{Q'} \left[ \frac{1}{p} \sum_{i=1}^{n_{1}} q_{i} T_{i} \right]$$
$$= \frac{\omega}{p} \sum_{i=1}^{n_{1}} q_{i} \sum_{j=1}^{i} y_{j}$$

substituting q<sub>i</sub> =

 $q_i = H(x_i) - H(x_{i-1})$ 

 $(\omega_1 + \omega)/\mu' = \frac{\omega}{p} \sum_{i=1}^{n_1} x_{i-1}^{-\alpha} y_i$ 

and  $F_{i}(x_{i}) = x_{i}^{-\alpha} = 1 - H(x_{i})$ ,

(ω<sub>1</sub>

we have

 $= \mu \sum_{1}^{n_{1}} x_{i-1}^{-\alpha} y_{i}$  (A1.5)

Substituting (A1.4) and (A1.5) into Equation (5) we obtain Equation (5').