

AN EXPERIMENTAL TIME-SHARING SYSTEM

Fernando J. Corbató, Marjorie Merwin-Daggett, Robert C. Daley

Computer Center, Massachusetts Institute of Technology
Cambridge, Massachusetts

Summary

It is the purpose of this paper to discuss briefly the need for time-sharing, some of the implementation problems, an experimental time-sharing system which has been developed for the contemporary IBM 7090, and finally a scheduling algorithm of one of us (FJC) that illustrates some of the techniques which may be employed to enhance and be analyzed for the performance limits of such a time-sharing system.

Introduction

The last dozen years of computer usage have seen great strides. In the early 1950's, the problems solved were largely in the construction and maintenance of hardware; in the mid-1950's, the usage languages were greatly improved with the advent of compilers; now in the early 1960's, we are in the midst of a third major modification to computer usage: the improvement of man-machine interaction by a process called time-sharing.

Much of the time-sharing philosophy, expressed in this paper, has been developed in conjunction with the work of an MIT preliminary study committee, chaired by H. Teager, which examined the long range computational needs of the Institute, and a subsequent MIT computer working committee, chaired by J. McCarthy. However, the views and conclusions expressed in this paper should be taken as solely those of the present authors.

Before proceeding further, it is best to give a more precise interpretation to time-sharing. One can mean using different parts of the hardware at the same time for different tasks, or one can mean several persons making use of the computer at the same time. The first meaning, often called multiprogramming, is oriented towards hardware efficiency in the sense of attempting to attain complete utilization of all components^{5,6,7,8}. The second meaning of time-sharing, which is meant here, is primarily concerned with the efficiency of persons trying to use a computer^{1,2,3,4}. Computer efficiency should still be considered but only in the perspective of the total system utility.

The motivation for time-shared computer usage arises out of the slow man-computer interaction rate presently possible with the bigger, more advanced computers. This rate has changed little (and has become worse in some cases) in the last decade of widespread computer use.¹⁰

In part, this effect has been due to the fact that as elementary problems become mastered on the computer, more complex problems immediately become of interest. As a result, larger and more complicated programs are written to take advantage of larger and faster computers. This process inevitably leads to more programming errors and a longer period of time required for debugging. Using current batch monitor techniques, as is done on most large computers, each program bug usually requires several hours to eliminate, if not a complete day. The only alternative presently available is for the programmer to attempt to debug directly at the computer, a process which is grossly wasteful of computer time and hampered seriously by the poor console communication usually available. Even if a typewriter is the console, there are usually lacking the sophisticated query and response programs which are vitally necessary to allow effective interaction. Thus, what is desired is to drastically increase the rate of interaction between the programmer and the computer without large economic loss and also to make each interaction more meaningful by extensive and complex system programming to assist in the man-computer communication.

To solve these interaction problems we would like to have a computer made simultaneously available to many users in a manner somewhat like a telephone exchange. Each user would be able to use a console at his own pace and without concern for the activity of others using the system. This console could as a minimum be merely a typewriter but more ideally would contain an incrementally modifiable self-sustaining display. In any case, data transmission requirements should be such that it would be no major obstacle to have remote installation from the computer proper.

1 The basic technique for a time-sharing system
 2 is to have many persons simultaneously using the
 3 computer through typewriter consoles with a
 4 time-sharing supervisor program sequentially
 5 running each user program in a short burst or
 6 quantum of computation. This sequence, which in
 7 the most straightforward case is a simple round-
 8 robin, should occur often enough so that each
 9 user program which is kept in the high-speed
 10 memory is run for a quantum at least once during
 11 each approximate human reaction time (~.2
 12 seconds). In this way, each user sees a computer
 13 fully responsive to even single key strokes each
 14 of which may require only trivial computation;
 15 in the non-trivial cases, the user sees a
 16 gradual reduction of the response time which is
 17 proportional to the complexity of the response
 18 calculation, the slowness of the computer, and
 19 the total number of active users. It should be
 20 clear, however, that if there are n users
 21 actively requesting service at one time, each
 22 user will only see on the average 1/n of the
 23 effective computer speed. During the period of
 24 high interaction rates while debugging programs,
 25 this should not be a hindrance since ordinarily
 26 the required amount of computation needed for
 27 each debugging computer response is small
 28 compared to the ultimate production need.

29 Not only would such a time-sharing system
 30 improve the ability to program in the conventional
 31 manner by one or two orders of magnitude, but
 32 there would be opened up several new forms of
 33 computer usage. There would be a gradual
 34 reformulation of many scientific and engineering
 35 applications so that programs containing decision
 36 trees which currently must be specified in
 37 advance would be eliminated and instead the
 38 particular decision branches would be specified
 39 only as needed. Another important area is that
 40 of teaching machines which, although frequently
 41 trivial computationally, could naturally
 42 exploit the consoles of a time-sharing system
 43 with the additional bonus that more elaborate
 44 and adaptive teaching programs could be used.
 45 Finally, as attested by the many small business
 46 computers, there are numerous applications in
 47 business and in industry where it would be
 48 advantageous to have powerful computing facilities
 49 available at isolated locations with only the
 50 incremental capital investment of each console.
 51 But it is important to realize that even without
 52 the above and other new applications, the major
 53 advance in programming intimacy available from
 54 time-sharing would be of immediate value to
 55 computer installations in universities, research
 56 laboratories, and engineering firms where
 57 program debugging is a major problem.

Implementation Problems

61 As indicated, a straightforward plan for
 62 time-sharing is to execute user programs for
 63 small quanta of computation without priority
 64 in a simple round-robin; the strategy of time-
 65 sharing can be more complex as will be shown
 66 later, but the above simple scheme is an
 67 adequate solution. There are still many
 68 problems, however, some best solved by hard-
 69 ware, others affecting the programming conven-
 70 tions and practices. A few of the more
 71 obvious problems are summarized:
 72

Hardware Problems:

73 1. Different user programs if simultan-
 74 eously in core memory may interfere with each
 75 other or the supervisor program so some form of
 76 memory protection mode should be available when
 77 operating user programs.
 78

79 2. The time-sharing supervisor may need
 80 at different times to run a particular program
 81 from several locations. (Loading relocation
 82 bits are no help since the supervisor does not
 83 know how to relocate the accumulator, etc.)
 84 Dynamic relocation of all memory accesses that
 85 pick up instructions or data words is one
 86 effective solution.
 87

88 3. Input-output equipment may be initiated
 89 by a user and read words in on another user
 90 program. A way to avoid this is to trap all
 91 input-output instructions issued by a user's
 92 program when operated in the memory protection
 93 mode.
 94

95 4. A large random-access back-up storage
 96 is desirable for general program storage files
 97 for all users. Present large capacity disc
 98 units appear to be adequate.
 99

100 5. The time-sharing supervisor must be
 101 able to interrupt a user's program after a
 102 quantum of computation. A program-initiated one-
 103 shot multivibrator which generates an interrupt
 104 at a fixed time later is adequate.
 105

106 6. Large core memories (e.g. a million
 107 words) would ease the system programming compli-
 108 cations immensely since the different active
 109 user programs as well as the frequently used
 110 system programs such as compilers, query programs,
 111 etc. could remain in core memory at all times.
 112

Programming Problems:

113 1. The supervisor program must do auto-
 114 matic user usage charge accounting. In general,
 115

116
 117
 118
 119
 120

1
2
3 the user should be charged on the basis of a
4 system usage formula or algorithm which should
5 include such factors as computation time, amount
6 of high-speed memory required, rent of secondary
7 memory storage, etc.

8 2. The supervisor program should coordinate
9 all user input-output since it is not desirable
10 to require a user program to remain constantly
11 in memory during input-output limited operations.
12 In addition, the supervisor must coordinate all
13 usage of the central, shared high-speed input-
14 output units serving all users as well as the
15 clocks, disc units, etc.

16 3. The system programs available must be
17 potent enough so that the user can think about
18 his problem and not be hampered by coding
19 details or typographical mistakes. Thus,
20 compilers, query programs, post-mortem programs,
21 loaders, and good editing programs are
22 essential.

23 4. As much as possible, the users should
24 be allowed the maximum programming flexibility
25 both in choices of language and in the absence
26 of restrictions.

27 Usage Problems

28 1. Too large a computation or excessive
29 typewriter output may be inadvertently requested
30 so that a special termination signal should be
31 available to the user.

32 2. Since real-time is not computer usage-
33 time, the supervisor must keep each user informed
34 so that he can use his judgment regarding loops,
35 etc.

36 3. Computer processor, memory and tape
37 malfunctions must be expected. Basic operational
38 questions such as "Which program is running?"
39 must be answerable and recovery procedures fully
40 anticipated.

41 An Experimental Time-Sharing System for the IBM 42 7090

43
44
45 Having briefly stated a desirable time-
46 sharing performance, it is pertinent to ask
47 what level of performance can be achieved with
48 existant equipment. To begin to answer this
49 question and to explore all the programming and
50 operational aspects, an experimental time-
51 sharing system has been developed. This system
52 was originally written for the IBM 709 but has
53 since been converted for use with the 7090
54 computer.
55
56
57
58
59
60

61
62 The 7090 of the MIT Computation Center has,
63 in addition to three channels with 19 tape units,
64 a fourth channel with the standard Direct Data
65 Connection. Attached to the Direct Data Connec-
66 tion is a real-time equipment buffer and control
67 rack designed and built under the direction of
68 H. Teager and his group. This rack has a variety
69 of devices attached but the only ones required
70 by the present systems are three flexowriter
71 typewriters. Also installed on the 7090 are two
72 special modifications (i.e. RPQ's): a standard
73 60 cycle accounting and interrupt clock, and a
74 special mode which allows memory protection,
75 dynamic relocation and trapping of all user
76 attempts to initiate input-output instructions.

77 In the present system the time-sharing
78 occurs between four users, three of whom are on-
79 line each at a typewriter in a foreground
80 system, and a fourth passive user of the back-
81 ground Fap-Mad-Madtran-BSS Monitor System similar
82 to the Fortran-Fap-BSS Monitor System (FMS) used
83 by most of the Center programmers and by many
84 other 7090 installations.

85 Significant design features of the fore-
86 ground system are:

87 1. It allows the user to develop programs
88 in languages compatible with the background
89 system,

90 2. Develop a private file of programs,
91
92 3. Start debugging sessions at the state
93 of the previous session, and
94

95 4. Set his own pace with little waste of
96 computer time.

97 Core storage is allocated such that all users
98 operate in the upper 27,000 words with the time-
99 sharing supervisor (TSS) permanently in the
100 lower 5,000 words. To avoid memory allocation
101 clashes, protect users from one another, and
102 simplify the initial 709 system organization,
103 only one user was kept in core memory at a
104 time. However, with the special memory protec-
105 tion and relocation feature of the 7090, more
106 sophisticated storage allocation procedures are
107 being implemented. In any case, user swaps are
108 minimized by using 2-channel overlapped magnetic
109 tape reading and writing of the pertinent
110 locations in the two user programs.

111 The foreground system is organized around
112 commands that each user can give on his type-
113 writer and the user's private program files
114 which presently (for want of a disc unit) are
115 kept on a separate magnetic tape for each user.
116

117 * This group is presently using another approach⁹
118 in developing a time-sharing system for the
119 MIT 7090.
120

1 For convenience the format of the private tape
 2 files is such that they are card images, have
 3 title cards with name and class designators and
 4 can be written or punched using the off-line
 5 equipment. (The latter feature also offers a
 6 crude form of large-scale input-output.) The
 7 magnetic tape requirements of the system are the
 8 seven tapes required for the normal functions of
 9 the background system, a system tape for the
 10 time-sharing supervisor that contains most of
 11 the command programs, and a private file tape
 12 and dump tape for each of the three foreground
 13 users.

14 The commands are typed by the user to the
 15 time-sharing supervisor (not to his own program)
 16 and thus can be initiated at any time regardless
 17 of the particular user program in memory. For
 18 similar coordination reasons, the supervisor
 19 handles all input-output of the foreground
 20 system typewriters. Commands are composed of
 21 segments separated by vertical strokes; the
 22 first segment is the command name and the
 23 remaining segments are parameters pertinent to
 24 the command. Each segment consists of the last
 25 6 characters typed (starting with an implicit
 26 6 blanks) so that spacing is an easy way to
 27 correct a typing mistake. A carriage return is
 28 the signal which initiates action on the command.
 29 Whenever a command is received by the supervisor,
 30 "WAIT", is typed back followed by "READY," when
 31 the command is completed. (The computer responses
 32 are always in the opposite color from the user's
 33 typing.) While typing, an incomplete command
 34 line may be ignored by the "quit" sequence of a
 35 code delete signal followed by a carriage return.
 36 Similarly after a command is initiated, it may
 37 be abandoned if a "quit" sequence is given. In
 38 addition, during unwanted command timeouts, the
 39 command and output may be terminated by pushing
 40 a special "stop output" button.

41 The use of the foreground system is initiated
 42 whenever a typewriter user completes a command
 43 line and is placed in a waiting command queue.
 44 Upon completion of each quantum, the time-sharing
 45 supervisor gives top priority to initiating any
 46 waiting commands. The system programs corres-
 47 ponding to most of the commands are kept on the
 48 special supervisor command system tape so that to
 49 avoid waste of computer time, the supervisor
 50 continues to operate the last user program until
 51 the desired command program on tape is positioned
 52 for reading. At this point, the last user is
 53 read out on his dump tape, the command program
 54 read in, placed in a working status and initiated
 55 as a new user program. However, before starting
 56 the new user for a quantum of computation, the
 57 supervisor again checks for any waiting command
 58 of another user and if necessary begins the look-
 59 ahead positioning of the command system tape
 60 while operating the new user.

Whenever the waiting command queue is
 empty, the supervisor proceeds to execute a
 simple round-robin of those foreground user
 programs in the working status queue. Finally,
 if both these queues are empty, the background
 user program is brought in and run a quantum at
 a time until further foreground system actively
 develops.

Foreground user programs leave the working
 status queue by two means. If the program
 proceeds to completion, it can reenter the
 supervisor in a way which eliminates itself and
 places the user in dead status; alternatively,
 by a different entry the program can be placed
 in a dormant status (or be manually placed by
 the user executing a quit sequence). The dormant
 status differs from the dead status in that the
 user may still restart or examine his program.

User input-output is through each type-
 writer, and even though the supervisor has a
 few lines of buffer space available, it is
 possible to become input-output limited.
 Consequently, there is an additional input-
 output wait status, similar to the dormant,
 which the user is automatically placed in by
 the supervisor program whenever input-output
 delays develop. When buffers become near
 empty on output or near full on input, the user
 program is automatically returned to the working
 status; thus waste of computer time is avoided.

Commands

To clarify the scope of the foreground
 system and to indicate the basic tools avail-
 able to the user, a list of the important
 commands follows along with brief summaries of
 their operations:

1. | α
 α = arbitrary text treated as a comment.
2. login | α | β
 α = user problem number
 β = user programmer number
 Should be given at beginning of each
 user's session. Rewinds user's private file tape;
 clears time accounting records.
3. logout
 Should be given at end of each user's
 session. Rewinds user's private file tape;
 punches on-line time accounting cards.
4. input
 Sets user in input mode and initiates
 automatic generation of line numbers. The user

1	types a card image per line according to a		
2	format appropriate for the programming language.		
3	(The supervisor collects these card images at		
4	the end of the user's private file tape.) When		
5	in the automatic input mode, the manual mode may		
6	be entered by giving an initial carriage return		
7	and typing the appropriate line number followed		
8	by and line for as many lines as desired. To		
9	reenter the automatic mode, an initial carriage		
10	return is given.		
11	The manual mode allows the user to over-		
12	write previous lines and to insert lines. (cf.		
13	File Command.)		
14			
15	5. edit α β		
16	α = title of file		
17	β = class of file		
18	The user is set in the automatic input		
19	mode with the designated file treated as initial		
20	input lines. The same conventions apply as to		
21	the input command.		
22			
23	6. file α β		
24	α = title to be given to file		
25	β = class of language used during input		
26	The created file will consist of the		
27	numbered input lines (i.e. those at the end of		
28	the user's private file tape) in <u>sequence</u> ; in		
29	the case of duplicate line numbers, the last		
30	version will be used. The line numbers will be		
31	written as sequence numbers in the corresponding		
32	card images of the file.		
33			
34	For convenience the following editing		
35	conventions apply to input lines:		
36			
37	a. an underline signifies the deletion of		
38	the previous characters of the line.		
39	b. a backspace signifies the deletion of		
40	the previous character in the field.		
41	The following formats apply:		
42			
43	a. FAP: symbol, tab, operation, tab,		
44	variable field and comment.		
45	b. MAD, MADTRAN, FORTRAN: statement label,		
46	tab, statement. To place a character in the		
47	continuation column: statement label, tab,		
48	backspace, character, statement.		
49	c. DATA: cols. 1-72.		
50			
51	7. fap α		
52	Causes the file designated as α , fap to		
53	be translated by the FAP translator (assembler).		
54	Files α , symtb and α ,bss are added to the user's		
55	private file tape giving the symbol table and		
56	the relocatable binary BSS form of the file.		
57			
58			
59			
60			
		8. mad α	61
		Causes file α ,mad to be translated by	62
		the MADtranslator (compiler). File α ,bss is	63
		created.	64
			65
		9. madtrn α	66
		Causes file α ,madtrn (i.e. a pseudo-	67
		Fortran language file) to be edited into an	68
		equivalent file α ,mad (added to the user's file)	69
		and translation occurs as if the command mad α	70
		had been given.	71
			72
		10. load α_1 α_2 ... α_n	73
		Causes the consecutive loading of files	74
		α_i ,bss (i=1,2,...,n). An exception occurs if α_1 =	75
		(libe), in which case file α_{i+1} ,bss is searched	76
		as a library file for all subprograms still	77
		missing. (There can be further library files.)	78
			79
		11. use α_1 α_2 ... α_n	80
		This command is used whenever a load or	81
		previous use command notifies the user of an	82
		incomplete set of subprograms. Same α_1 conven-	83
		tions as for load.	84
			85
		12. start α β	86
		Starts the program setup by the load	87
		and use commands (or a dormant program) after	88
		first positioning the user private file tape in	89
		front of the title card for file α , β . (If β is	90
		not given, a class of data is assumed; if both	91
		α and β are not given, no tape movement occurs	92
		and the program is started.)	93
			94
		13. pm α	95
		α = "lights", "stomap", or the usual	96
		format of the standard Center post-mortem (F2PM)	97
		request: subprogram name loc ₁ loc ₂ mode	98
		direction where mode and direction are optional.	99
		Produces post-mortem of user's dormant	100
		program according to request specified by α .	101
		(E.g. matrix 5 209 flo rev will cause to	102
		be printed on the user's typewriter the contents	103
		of subprogram "matrix" from relative locations	104
		5 to 209 in floating point form and in reverse	105
		sequence.)	106
			107
		14. skipppm	108
		Used if a pm command is "quit" during	109
		output and the previous program interruption is to	110
		be restarted.	111
			112
		15. listf	113
		Types out list of all file titles on	114
		user's private file tape.	115
			116
			117
			118
			119
			120

1 16. printf | α | β | γ
2 Types out file α, β starting at line
3 number γ . If γ is omitted, the initial line is
4 assumed. Whenever the user's output buffer fills,
5 the command program goes into an I/O wait status
6 allowing other users to time-share until the
7 buffer needs refilling.
8
9 17. xdump | α | β
10 Creates file α, β (if β omitted, xdump
11 assumed) on user's private file tape consisting
12 of the complete state of the user's last dormant
13 program.
14
15 18. xdump | α | β
16 Inverse of xdump command in that it
17 resets file α, β as the user's program, starting
18 it where it last left off.
19
20 Although experience with the system to date
21 is quite limited, first indications are that
22 programmers would readily use such a system if it
23 were generally available. It is useful to ask,
24 now that there is some operating experience with
25 the 7090 system, what observations can be made.
26 An immediate comment is that once a user gets
27 accustomed to computer response, delays of even
28 a fraction of a minute are exasperatingly long,
29 an effect analogous to conversing with a slow-
30 speaking person. Similarly, the requirement that
31 a complete typewritten line rather than each
32 character be the minimum unit of man-computer
33 communication is an inhibiting factor in the
34 sense that a press-to-talk radio-telephone con-
35 versation is more stilted than that of an
36 ordinary telephone. Since maintaining a rapid
37 computer response on a character by character
38 basis requires at least a vestigial response
39 program in core memory at all times, the straight-
40 forward solution within the present system is to
41 have more core memory available. At the very
42 least, an extra bank of memory for the time-
43 sharing supervisor would ease compatibility prob-
44 lems with programs already written for 32,000
45 word 7090's.
46
47 For reasons of expediency, the weakest
48 portions of the present system are the conventions
49 for input, editing of user files, and the degree
50 of rapid interaction and intimacy possible while
51 debugging. Since to a large extent these areas
52 involve the taste, habits, and psychology of the
53 users, it is felt that proper solutions will
54 require considerable experimentation and prag-
55 matic evaluation; it is also clear that these
56 areas cannot be treated in the abstract for the
57 programming languages used will influence greatly
58 the appropriate techniques. A greater use of
59 symbolic referencing for locations, program names
60 and variables is certainly desired; symbolic post-
mortem programs, trace programs, and before-and-
after differential dump programs should play
useful roles in the debugging procedures.

In the design of the present system, great
care went into making each user independent of
the other users. However, it would be a useful
extension of the system if this were not always
the case. In particular, when several consoles
are used in a computer controlled group such as
in management or war games, in group behavior
studies, or possibly in teaching machines, it
would be desirable to have all the consoles
communicating with a single program.

Another area for further improvement within
the present system is that of file maintenance,
since the presently used tape units are a hind-
rance to the easy deletion of user program files.
Disc units will be of help in this area as well
as with the problem of consolidating and
scheduling large-scale central input-output
generated by the many console users.

Finally, it is felt that it would be desir-
able to have the distinction between the fore-
ground and background systems eliminated. The
present-day computer operator would assume the
role of a stand-in for the background users,
using an operator console much like the other
user consoles in the system, mounting and de-
mounting magnetic tapes as requested by the
supervisor, receiving instructions to read card
decks into the central disc unit, etc. Similarly
the foreground user, when satisfied with his
program, would by means of his console and the
supervisor program enter his program into the
queue of production background work to be
performed. With these procedures implemented
the distinction of whether one is time-sharing
or not would vanish and the computer user would
be free to choose in an interchangeable way that
mode of operation which he found more suitable
at a particular time.

A Multi-Level Scheduling Algorithm

Regardless of whether one has a million
word core memory or a 32,000 word memory as
currently exists on the 7090, one is inevitably
faced with the problem of system saturation
where the total size of active user programs
exceeds that of the high-speed memory or there
are too many active user programs to maintain
an adequate response at each user console.
These conditions can easily arise with even a
few users if some of the user programs are
excessive in size or in time requirements. The
predicament can be alleviated if it is assumed
that a good design for the system is to have a
saturation procedure which gives graceful de-
gradation of the response time and effective
real-time computation speed of the large and
long-running users.

To show the general problem, Figure 1 qualitatively gives the user service as a function of n, the number of active users. This service parameter might be either of the two key factors: computer response time or n times the real-time computation speed. In either case there is some critical number of active users, N, representing the effective user capacity, which causes saturation. If the strategy near saturation is to execute the simple round-robin of all users, then there is an abrupt collapse of service due to the sudden onset of the large amount of time required to swap programs in-and-out of the secondary memory such as a disc or drum unit. Of course, Figure 1 is quite qualitative since it depends critically on the spectrum of user program sizes as well as the spectrum of user operating times.

To illustrate the strategy that can be employed to improve the saturation performance of a time-sharing system, a multi-level scheduling algorithm is presented. This algorithm also can be analyzed to give broad bounds on the system performance.

The basis of the multi-level scheduling algorithm is to assign each user program as it enters the system to be run (or completes a response to a user) to an lth level priority queue. Programs are initially entered into a level l₀, corresponding to their size such that

$$l_0 = \left\lceil \log_2 \left(\left\lceil \frac{w_p}{w_q} \right\rceil + 1 \right) \right\rceil \quad (1)$$

where w_p is the number of words in the program, w_q is the number of words which can be transmitted in and out of the high-speed memory from the secondary memory in the time of one quantum, q, and the bracket indicates "the integral part of". Ordinarily the time of a quantum, being the basic time unit, should be as small as possible without excessive overhead losses when the supervisor switches from one program in high-speed memory to another. The process starts with the time-sharing supervisor operating the program at the head of the lowest level occupied queue, l, for up to 2^l quanta of time and then if the program is not completed (i.e. has not made a response to the user) placing it at the end of the l+1 level queue. If there are no programs entering the system at levels lower than l, this process proceeds until the queue at level l is exhausted; the process is then iteratively begun again at level l+1, where now each program is run for 2^{l+1} quanta of time. If during the execution of the 2^l quanta of a program at level l, a lower level, l', becomes occupied, the current user is replaced at the head of the lth queue and the process is reinitiated at level l'.

Similarly, if a program of size w_p at level l, during operation requests a change in memory size from the time-sharing supervisor, then the enlarged (or reduced) version of the program should be placed at the end of the l'' queue where

$$l'' = l + \left\lceil \log_2 \left(\left\lceil \frac{w_p''}{w_p} \right\rceil + 1 \right) \right\rceil \quad (2)$$

Again the process is re-initiated with the head-of-the-queue user at the lowest occupied level of l'.

Several important conclusions can be drawn from the above algorithm which allow the performance of the system to be bounded.

Computational Efficiency

1. Because a program is always operated for a time greater than or equal to the swap time (i.e. the time required to move the program in and out of secondary memory), it follows that the computational efficiency never falls below one-half. (Clearly, this fraction is adjustable in the formula for the initial level, l₀.) An alternative way of viewing this bound is to say that the real-time computing speed available to one out of n active users is no worse than if there were 2n active users all of whose programs were in the high-speed memory.

Response Time

2. If the maximum number of active users is N, then an individual user of a given program size can be guaranteed a response time,

$$t_r \leq 2Nq \left(\left\lceil \frac{w_p}{w_q} \right\rceil + 1 \right) \quad (3)$$

since the worst case occurs when all competing user programs are at the same level. Conversely, if t_r is a guaranteed response of arbitrary value and the largest size of program is assumed, then the maximum permissible number of active users is bounded.

Long Runs

3. The relative swap time on long runs can be made vanishingly small. This conclusion follows since the longer a program is run, the higher the level number it cascades to with a correspondingly smaller relative swap time. It is an important feature of the algorithm that long runs must in effect prove they are long so that programs which have an unexpected demise are detected quickly. In order that there be a finite number of levels, a maximum level number, L, can be established such that the asymptotic swap overhead is some arbitrarily small percentage, p:

Note: formula (2) is incorrect. To correct, remove "+ 1" and remove inner brackets around ratio of program sizes.

$$L = \left[\log_2 \left(\left[\frac{w_{pmax}}{pw_q} \right] + 1 \right) \right] \quad (4)$$

where w_{pmax} is the size of the largest possible program.

Multi-level vs. Single-level Response Times

4. The response time for programs of equal size, entering the system at the same time, and being run for multiple quanta, is no worse than approximately twice the response-time occurring in a single quanta round-robin procedure. If there are n equal sized programs started in a queue at level ℓ , then the worst case is that of the end-of-the-queue program which is ready to respond at the very first quantum run at the $\ell+j$ level. Using the multi-level algorithm, the total delay for the end-of-the-queue program is by virtue of the geometric series of quanta:

$$T_m \sim q2^\ell \{n(2^j-1) + (n-1)2^j\} \quad (5)$$

Since the end-of-the-queue user has computed for a time of $2^j(2^j-1)$ quanta, the equivalent single-level round-robin delay before a response is:

$$T_s \sim q2^\ell \{n(2^j-1)\} \quad (6)$$

Hence

$$\frac{T_m}{T_s} \sim 1 + \left(\frac{n-1}{n}\right) \left(\frac{2^j}{2^j-1}\right) \sim 2 \quad (7)$$

and the assertion is shown. It should be noted that the above conditions, where program swap times are omitted, which are pertinent when all programs remain in high-speed memory, are the least favorable for the multi-level algorithm; if swap times are included in the above analysis, the ratio of T_m/T_s can only become smaller and may become much less than unity. By a similar analysis it is easy to show that even in the unfavorable case where there are no program swaps, head-of-the-queue programs that terminate just as the $2^{\ell+j}$ quanta are completed receive under the multi-level algorithm a response which is twice as fast as that under the single-level round-robin (i.e. $T_m/T_s = 1/2$).

Highest Serviced Level

5. In the multi-level algorithm the level classification procedure for programs is entirely automatic, depending on performance and program size rather than on the declarations (or hopes) of each user. As a user taxes the system, the degradation of service occurs progressively starting with the higher level users of either large or long-running programs; however, at some level no user programs may be run because of too many active users at lower levels. To determine

a bound on this cut-off point we consider N active users at level ℓ each running 2^ℓ quanta, terminating, and reentering the system again at level ℓ at a user response time, t_u , later. If there is to be no service at level $\ell+1$, then the computing time, $Nq2^\ell$, must be greater than or equal to t_u . Thus the guaranteed active levels, ℓ_a , are given by the relation:

$$\ell_a \leq \left[\log_2 \left(\frac{t_u}{Nq} \right) \right] \quad (8)$$

In the limit, t_u could be as small as a minimum user reaction time ($\sim .2$ sec.), but the expected value would be several orders of magnitude greater as a result of the statistics of a large number of users.

The multi-level algorithm as formulated above makes no explicit consideration of the seek or latency time required before transmission of programs to and from disc or drum units when they are used as the secondary memory, (although formally the factor w_q could contain an average figure for these times). One simple modification to the algorithm which usually avoids wasting the seek or latency time is to continue to operate the last user program for as many quanta as are required to ready the swap of the new user with the least priority user; since ordinarily only the higher level number programs would be forced out into the secondary memory, the extended quanta of operation of the old user while seeking the new user should be but a minor distortion of the basic algorithm.

Further complexities are possible when the hardware is appropriate. In computers with input-output channels and low transmission rates to and from secondary memory, it is possible to overlap the reading and writing of the new and old users in and out of high-speed memory while operating the current user. The effect is equivalent to using a drum giving 100% multiplexor usage but there are two liabilities, namely, no individual user can utilize all the available user memory space and the look-ahead procedure breaks down whenever an unanticipated scheduling change occurs (e.g. a program terminates or a higher-priority user program is initiated).

Complexity is also possible in storage allocation but certainly an elementary procedure and a desirable one with a low-transmission rate secondary memory is to consolidate in a single block all high-priority user programs whenever sufficient fragmentary unused memory space is available to read in a new user program. Such a procedure is indicated in the flow diagram of the multi-level scheduling algorithm which is given as Figure 2.

1 It should also be noted that Figure 2 only
 2 accounts for the scheduling of programs in a
 3 working status and still does not take into
 4 account the storage allocation of programs which
 5 are in a dormant (or input-output wait status).
 6 One systematic method of handling this case is
 7 to modify the scheduling algorithm so that
 8 programs which become dormant at level l are
 9 entered into the queue at level $l+1$. The
 10 scheduling algorithm proceeds as before with the
 11 dormant programs continuing to cascade but not
 12 operating when they reached the head of a queue.
 13 Whenever a program must be removed from high-
 14 speed memory, a program is selected from the end-
 15 of-the-queue of the highest occupied level
 16 number.

17 Finally, it is illuminating to apply the
 18 multi-level scheduling algorithm bounds to the
 19 contemporary IBM 7090. The following approximate
 20 values are obtained:

21 $q = 16$ m.s. (based on 1% switching overhead)

22 $w_q = 120$ words (based on one IBM 1301 model
 23 2 disc unit without seek or latency
 24 times included)

25 $t_r \leq 8Nf_{\text{sec.}}$ (based on programs of (32k)f
 26 words)

27 $l_a \leq \log_2 (1000/N)$ (based on $t_u = 16$ sec.)

28 $l_o \leq 8$ (based on a maximum program size of
 29 32K words)

30
 31 Using the arbitrary criteria that programs
 32 up to the maximum size of 32,000 words should
 33 always get some service, which is to say that
 34 $\max l_a = \max l_o$, we deduce as a conservative
 35 estimate that N can be 4 and that at worst the
 36 response time for a trivial reply will be 32
 37 seconds.
 38

39 The small value of N arrived at is a direct
 40 consequence of the small value of w_q that results
 41 from the slow disc word transmission rate. This
 42 rate is only 3.3% of the maximum core memory
 43 multiplexor rate. It is of interest that using
 44 high-capacity high-speed drums of current design
 45 such as in the Sage System or in the IBM Sabre
 46 System it would be possible to attain nearly
 47 100% multiplexor utilization and thus multiply
 48 w_q by a factor of 30. It immediately follows
 49 that user response times equivalent to those
 50 given above with the disc unit would be given
 51 to 30 times as many persons or to 120 users; the
 52 total computational capacity, however, would not
 53 change.
 54
 55
 56
 57
 58
 59
 60

61 In any case, considerable caution should be
 62 used with capacity and computer response time
 63 estimates since they are critically dependent
 64 upon the distribution functions for the user
 65 response time, t_u , and the user program size,
 66 w_p , and the computational capacity requested by
 67 each user. Past experience using conventional
 68 programming systems is of little assistance be-
 69 cause these distribution functions will depend
 70 very strongly upon the programming systems made
 71 available to the time-sharing users as well as
 72 upon the user habit patterns which will gradually
 73 evolve.
 74

75 Conclusions

76 In conclusion, it is clear that contemporary
 77 computers and hardware are sufficient to allow
 78 moderate performance time-sharing for a limited
 79 number of users. There are several problems
 80 which can be solved by careful hardware design,
 81 but there are also a large number of intricate
 82 system programs which must be written before one
 83 has an adequate time-sharing system. An import-
 84 ant aspect of any future time-shared computer is
 85 that until the system programming is completed,
 86 especially the critical time-sharing supervisor,
 87 the computer is completely worthless. Thus, it
 88 is essential for future system design and imple-
 89 mentation that all aspects of time-sharing system
 90 problems be explored and understood in prototype
 91 form on present computers so that major advances
 92 in computer organization and usage can be made.
 93

94 Acknowledgements

95 The authors wish to thank Bernard Galler,
 96 Robert Graham and Bruce Arden, of the University
 97 of Michigan, for making the MAD compiler available
 98 and for their advice with regard to its adaption
 99 into the present time-sharing system. The
 100 version of the Madtran Fortran-to-Mad editor
 101 program was generously supplied by Robert Rosin
 102 of the University of Michigan. Of the MIT
 103 Computation Center staff, Robert Creasy was of
 104 assistance in the evaluation of time-sharing
 105 performance, Lynda Korn is to be credited for her
 106 contributions to the pm and madtran commands, and
 107 Evelyn Dow for her work on the fap command.
 108

109 References

- 110
 111 1. Strachey, C., "Time Sharing in Large
 112 Fast Computers," Proceedings of the International
 113 Conference on Information Processing, UNESCO
 114 (June, 1959), Paper B.2.19.
 115
 116
 117
 118
 119
 120

1 2. Licklider, J. C. R., "Man-Computer
 2 Symbiosis," IRE Transactions on Human Factors in
 3 Electronics, HFE-1, No. 1 (March, 1960), 4-11.
 4
 5 3. Brown, G., Licklider, J. C. R.,
 6 McCarthy, J., and Perlis, A., lectures given
 7 spring, 1961, Management and the Computer of the
 8 Future, (to be published by the M.I.T. Press,
 9 March, 1962).
 10 4. Corbato, F. J., "An Experimental Time-
 11 Sharing System," Proceedings of the IBM University
 12 Director's Conference, July, 1961 (to be pub-
 13 lished).
 14 5. Schmitt, W. F., Tonik, A. B., "Sympa-
 15 thetically Programmed Computers," Proceedings
 16 of the International Conference on Information
 17 Processing, UNESCO, (June, 1959) Paper B.2.18.

61 6. Codd, E. F., "Multiprogram Scheduling,"
 62 Communications of the ACM, 3, 6 (June, 1960),
 63 347-350.
 64
 65 7. Heller, J., "Sequencing Aspects of
 66 Multiprogramming," Journal of the ACM, 8, 3
 67 (July, 1961), 426-439.
 68
 69 8. Leeds, H. D., Weinberg, G. M., "Multi-
 70 programming," Computer Programming Fundamentals,
 71 356-359, McGraw-Hill (1961).
 72
 73 9. Teager, H. M., "Real-Time Time-Shared
 74 Computer Project," Communications of the ACM, 5,
 75 1 (January, 1962) Research Summaries, 62.
 76
 77 10. Teager, H. M., McCarthy, J., "Time-
 78 Shared Program Testing," paper delivered at the
 79 14th National Meeting of the ACM (not published).
 80
 81
 82
 83
 84
 85
 86
 87
 88
 89
 90
 91
 92
 93
 94
 95
 96
 97
 98
 99
 100
 101
 102
 103
 104
 105
 106
 107
 108
 109
 110
 111
 112
 113
 114
 115
 116
 117
 118
 119
 120

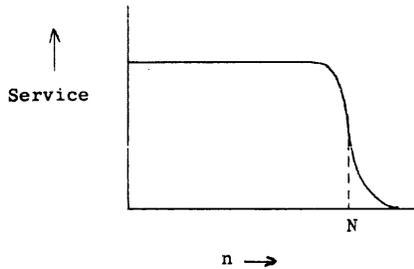


Figure 1. Service vs. Number of Active Users

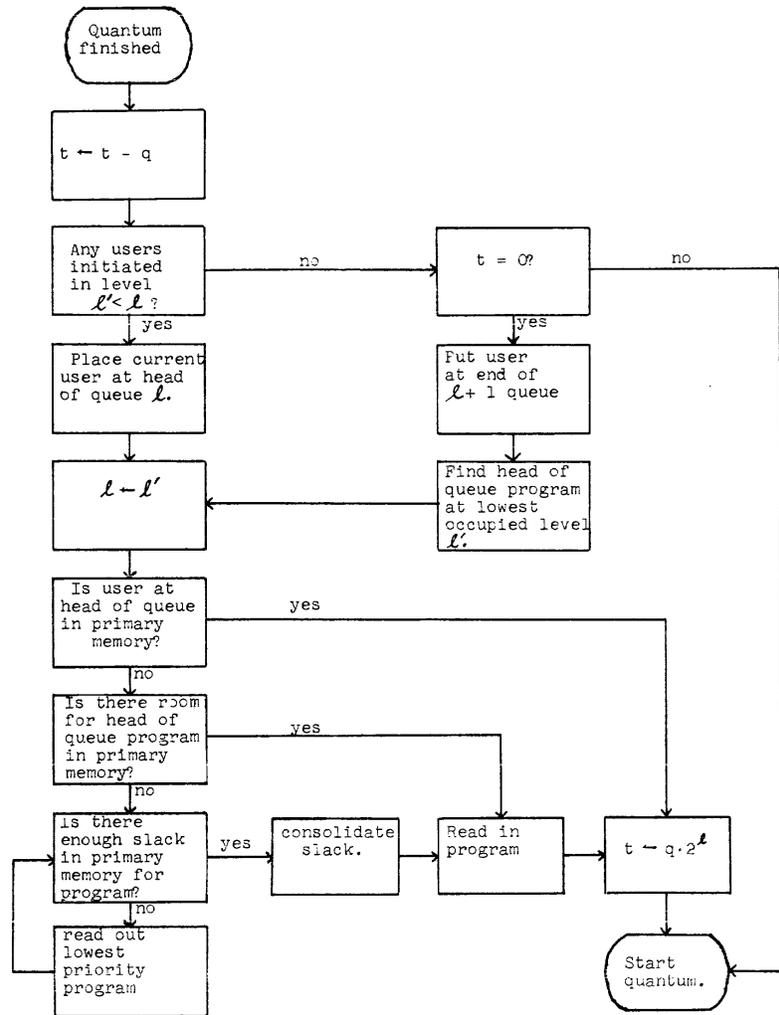


Figure 2. Flow Chart of Multi-Level Scheduling Algorithm