

TENEX, a Paged Time Sharing System for the PDP-10

Daniel G. Bobrow, Jerry D. Burchfiel,
Daniel L. Murphy, and Raymond S. Tomlinson
Bolt Beranek and Newman Inc.*

19 TENEX is a new time sharing system implemented
20 on a DEC PDP-10 augmented by special paging
21 hardware developed at BBN. This report specifies a
22 set of goals which are important for any time sharing
23 system. It describes how the TENEX design and
24 implementation achieve these goals. These include
25 specifications for a powerful multiprocess large
26 memory virtual machine, intimate terminal interaction,
27 comprehensive uniform file and I/O capabilities, and
28 clean flexible system structure. Although the
29 implementation described here required some
30 compromise to achieve a system operational within
31 six months of hardware checkout, TENEX has met its
32 major goals and provided reliable service at several sites
33 and through the ARPA network.

34 **Key Words and Phrases:** TENEX, paging, virtual
35 machines, time sharing system, scheduling algorithm,
36 process structure, PDP-10

37 **CR Categories:** 2.44, 4.32, 4.39, 4.42

47 Copyright © 1972, Association for Computing Machinery, Inc.
48 General permission to republish, but not for profit, all or part
49 of this material is granted, provided that reference is made to this
50 publication, to its date of issue, and to the fact that reprinting
51 privileges were granted by permission of the Association for Computing
52 Machinery.

52 Presented at the Third Annual Symposium on Operating Systems Principles, Palo Alto, California, October 18-20, 1971.

1. Introduction

TENEX is a new time sharing operating system implemented on the DEC PDP-10. It was developed because no existing system of the appropriate size and cost could meet the requirements of the research projects at BBN. During the development phase of TENEX we formulated a set of goals which we hoped would produce a workable and well-integrated system, as well as help us to realize our specific requirements. Our background included knowledge and use of a number of other systems including the DEC PDP-1 systems designed at BBN[1], the Berkeley system for the SDS 940[10], MIT CTSS[3], the DEC 10/50 system[5], and MULTICS[4]. We used good design ideas from each of these systems, and tried to avoid what we felt were problems of operation and implementation which we saw in these systems. The scale of the system and available resources imposed certain constraints on the implementation, including (1) that minimal change be made to the PDP-10 processor and none to the basic address computation, and (2) that the system had to be in service for users within six months of the operation of the hardware and at that time dominate our previous service, the DEC 10/50 and the SDS 940 systems.

Our design goals, presented below, are generally understood. Some, however, are often overlooked and usually not emphasized. We considered all of these important in the design of TENEX. Our design goals

* Computer Science Division, 50 Moulton Street, Cambridge, MA 02138. The work reported here was supported in part by the Advanced Research Projects Agency of the DOD, and in part by BBN.

1 fall into three broad categories with several specific
2 objectives under each.

3 I. State of the art virtual machine.

4 a. Paged virtual address space equal to or greater
5 than the addressing capability of the processor with
6 full provision for protection and sharing.

7 b. Multiple process capability in virtual machine
8 with appropriate communication facilities.

9 c. File system integrated into virtual address space,
10 built on multilevel symbolic directory structure with
11 protection, and providing consistent access to all
12 external I/O devices and data streams.

13 d. Extended instruction repertoire making available
14 many common operations as single instructions.

15 II. Good human engineering throughout system.

16 a. An executive command language interpreter which
17 provides direct access to a large variety of small,
18 commonly used system functions, and access to and
19 control over all other subsystems and user programs.
20 Command language forms should be extremely ver-
21 satile, adapting to the skill and experience of the user.

22 b. Terminal interface design which facilitates intimate
23 interaction between program and user, provides ex-
24 tensive interrupt capability, and full ASCII character set.

25 c. Virtual machine functions which provide all neces-
26 sary options, with reasonable default values simplify-
27 ing common cases, and require no system-created
28 objects to be placed in the user address space.

29 d. The system should encourage and facilitate co-
30 operation among users as well as provide protection
31 against undesired interaction.

32 III. The system must be implementable, maintain-
33 able, and modifiable.

34 a. Software must be modular with well-defined inter-
35 faces and with provision for adding or changing mod-
36 ules clearly considered.

37 b. Software must be debuggable and reliable, allow-
38 ing use of available debugging aids and including in-
39 ternal redundancy checks.

40 c. System should run efficiently, allow dynamic
41 manual adjustment of service if desired, and allow
42 extensive reconfiguration without reassembly.

43 d. System should contain instrumentation to clearly
44 indicate performance.

45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

2. Hardware Development for TENEX

Hardware development and modification for TENEX was limited to that necessary to achieve the goals specified above. This effort included an address mapping (paging) device implemented with then current DEC modules and some changes to the PDP-10 processor. A hard limit on the latter resulted from the lack of physical room available in the processor. Both processor changes and paging facilities had to be designed to allow the standard DEC time sharing software and diagnostics to run.

2.1 The BBN Pager

The BBN pager is an interface between the PDP-10 processor and the memory bus. It provides individual mapping (relocation) of each page (512 words) of both user and monitor address spaces using separate maps for each. The pager uses "associative registers" and core memory tables to store the mapping information. On each memory request from the processor, the 9 high-order bits of the address and the request-type level (read, write, execute) are compared in parallel with the contents of each associative register. If a match is found, the register containing the match also contains 11 high-order address bits to reference up to 1,048,576 words of physical core.

If no match is found, reference is made to a 512 word "page table" in physical core memory. The word selected in this page table is determined by a dispatch based on the original 9 high-order address bits. In the simple case of a private page which is in core, the 11 high-order address bits and protection bits are found in this word and are automatically loaded into an associative register by the pager.

There are three other cases:

1. The page is not in core, is protected from the requested type of access, or is nonexistent; in this case a page fault (trap) will occur.
2. The page is shared; in this case the map contains a "shared" pointer to a system table which contains the location information for the page.
3. The page belongs to another process; in this case, the entry contains an "indirect" pointer to an entry in another page table from which the location information is obtained.

The goal of program (code and data) sharing was given extensive consideration in the design of the BBN pager. The indirect and shared pointer mechanism allows pages to be actively shared (be represented in more than one address space) but still have the current address (core or secondary storage) stored in only one place. This allows memory control tables and data structures to be kept simple and the memory management software to move pages without extensive computational overhead. The pager permits individual pages to be shared for write as well as read references, so two or more processes may communicate by sharing

61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
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

1 a common page into which any or all may write.
2 Rather than enforce a discipline of pure procedures,
3 with private data in another segment, a unique "copy-
4 on-write" facility allows users to share large portions
5 of an address space containing procedures and data,
6 and to obtain private copies of only those pages which
7 are changed. This is implemented through an inde-
8 pendent per-page status bit, available to users, which
9 will produce a trap on a write reference, and a system
10 procedure by which a private copy is then created.
11 For example, this permits shared programs to be
12 prepared with preconstructed data areas, which will
13 be kept shared if not modified and will be put in private
14 storage *if changed*.

15 One final, unique feature is that the pager maintains
16 a record of the activity of the pages in core memory
17 in a "core status" table. The pager notes when a page
18 has been referenced, which processes have used that
19 page, and whether the page has been written into.
20 This information is used by the memory management
21 software to be described.

22 2.2 Processor Modifications

23 Except for the pager trapping facilities, all of the
24 TENEX virtual machine facilities (monitor calls) are
25 reached via a new system call instruction, JSYS, added
26 to the PDP-10 processor. JSYS provides an independent
27 transfer mechanism into the monitor which does not
28 conflict with "UUO" system calls used by DEC software.
29 It accomplishes a transfer from the user program to the
30 specified monitor routine in one instruction time via a
31 block of cells called the JSYS transfer vector. The JSYS
32 address provides the index to the proper transfer vector
33 entry. The state of the processor, including the return
34 address, is stored in a location specified by the transfer
35 vector, usually in a separate data area, so that a JSYS
36 call is suitable for reentrant code. The JSYS transfer
37 vector occupies exactly one page in the monitor space
38 and could be mapped independently for each process,
39 but this is not done in the current system.

40 There exists a context, either user or monitor, for
41 each instruction execution. The JSYS system call may
42 be executed in either context, with the "callee" operat-
43 ing in monitor context. One hardware modification
44 to the PDP-10 adds a bit to the state word of the proces-
45 sor to record the context of the "caller." Two ways of
46 accessing the caller memory context were added to the
47 PDP-10. One is an execute instruction which allows
48 current or previous context to be specified for each
49 memory reference of the object instruction. The second
50 access instruction group moves data between the
51 AC's (general registers) of the current context and
52 memory of the previous context. AC's from the pre-
53 vious context are accessed using a pager " AC base reg-
54 ister" which specifies the location of the stored AC's.
55
56
57
58
59
60

3. The TENEX Virtual Machine

61
62
63 A user process running under TENEX operates on a
64 virtual machine similar to a PDP-10 arithmetic processor
65 with 256K of virtual memory. The paging hardware
66 traps processor references to any data not in core, and a
67 core manager performs the necessary I/O to make the
68 referenced page available. Such traps are invisible to
69 the user process.

70 The virtual processor does not make available to
71 the user the direct I/O instructions of the PDP-10.
72 But through instructions which call monitor routines,
73 the virtual machine provides facilities that are con-
74 siderably more powerful and sophisticated than typical
75 hardware configurations used directly.

3.1 Virtual Memory Structure

76
77 The TENEX virtual memory may be viewed as a
78 linear address space of 256K words, and programs may
79 use it in this fashion. However, the existence of the
80 paging hardware means that the monitor must deal
81 with memory in pages of 512 words, and some of the
82 power which the mapping hardware provides is acces-
83 sible to a user program.

84 The contents of the virtual memory are specified by
85 the *virtual memory map* of 512 slots which the user may
86 read or write via monitor calls. The contents of each slot
87 specify the page in that position in the virtual address
88 space, and the type of access allowable (read and/or
89 write and/or execute) for that page. In the simplest
90 case, a map slot may contain (a pointer to) a private
91 page, i.e. a page shared with no other processes in
92 the system. A private page is automatically created
93 whenever a process makes a reference to a page for
94 which the map slot is empty. A slot may also con-
95 tain an indirect pointer to a page in this or some other
96 process. A memory reference to a location in such a
97 page will be executed just as though the instruction
98 had directly addressed the page pointed to. Any change
99 made to the page by either process will be seen by
100 both processes. If the owner of the page changes the
101 contents of his memory map, then the process with
102 the indirect pointer will see the change. A virtual
103 memory slot may also contain a pointer to a page
104 from a file in the file system, as discussed later.

3.2 Job Structure

105
106 A job is a set of one or more hierarchically related
107 processes, and it has the following attributes.

- 108 1. The name of user who initiated the job.
- 109 2. An account number to charge costs associated with
110 use of system resources.
- 111 3. Some open files.
- 112 4. A hierarchy of running and/or suspended processes.

113 A job may also have one or more terminal or
114 other devices assigned and attached.

115
116 **3.2.1 Process Hierarchy.** TENEX permits each job
117 to have multiple simultaneously runnable processes.
118

1 The relationships among them are defined by a struc-
2 ture which looks like an inverted tree defined by the
3 capability for direct control and killing. A process
4 always has exactly one immediately superior process
5 and may have one or more inferior processes. Two
6 processes are said to be parallel if they have the same
7 immediate superior. In TENEX, a process may create
8 processes inferior, but *not* parallel or superior in the
9 structure. A process can communicate with other
10 members of the structure by (a) sharing memory,
11 (b) direct control (superior to inferior only), or (c)
12 pseudo (software simulated) interrupts as described
13 in Section 3.3.

14 Although not completely general, a tree structure
15 process heirarchy implicitly provides the protection
16 and reference facilities that are wanted in most appli-
17 cations. These include referencing inferior processes
18 as a class for freezing, killing, and resuming; fielding
19 of interrupts and special conditions by a superior
20 process; and protection of the superior process from
21 inferiors.

22 Currently in TENEX, multiple processes are used:
23 1. To enable the EXEC to run user programs, handling
24 faults, and servicing user requested interrupts.
25 2. To allow programs to block for an arbitrary set of
26 events; one process waits for each event and signals
27 the main process when it occurs.
28 3. To implement an invisible debugging program,
29 completely protected from malfunction of the program
30 under test.

31 3.3 Pseudo Interrupt System

32 TENEX provides a facility for a process to receive
33 asynchronous signals from other processes or from
34 terminals or as the result of its own execution. The
35 various processes in a job may explicitly direct inter-
36 rupts to each other for purposes of communication.
37 A process may enable an interrupt which will occur
38 whenever the user hits a particular key on the control-
39 ling terminal. Finally, a process may use the pseudo
40 interrupt system to detect any of a set of unusual con-
41 ditions, including illegal references to memory, proc-
42 essor overflow conditions, end-of-file, and data errors.

44 3.4 Other Monitor Functions

45 Other functions which form a part of the virtual
46 machine include:
47 1. Functions which provide information to the pro-
48 gram about the state of the system or job (time of day,
49 runtime used, name of user, etc.).
50 2. Functions which save and restore the computational
51 environment of a process to allow restarting of a sus-
52 pended program.
53 3. Functions which provide frequently needed forms of
54 I/O conversions, such as fixed or floating point number
55 input and output, and date and time to string conver-
56 sions.
57

3.5 Backward Compatibility (DEC 10/50 Monitors) 61

62 Since TENEX was being implemented on a machine
63 for which a large useful program library existed, mostly
64 for use under the DEC 10/50 time sharing monitor,
65 we felt it was highly desirable to be able to run such
66 programs under the new monitor system. We felt it
67 should be possible to run *binary* images of old pro-
68 grams, i.e. without reassembling.

69 Toward this end, all of the TENEX monitor calls
70 were implemented with the JSYS instruction, reserving all
71 old monitor calls for their previous use. Secondly,
72 routines were designed which implemented all of the
73 existing 10/50 monitor calls in terms of the available
74 TENEX monitor calls. This set of routines implements
75 all of the functions available in the 10/50 monitor
76 except those specifically intended for the maintenance
77 of the system. Assembled together as a compatability
78 package, they occupy slightly less than 2.5K of core.
79 The package is kept as a core image file and is never
80 seen by programs which use only TENEX monitor calls.
81 However, the functions are automatically made avail-
82 able to 10/50 type programs by the monitor. When a
83 program makes its first 10/50 type monitor call, the
84 TENEX monitor maps the compatability package into a
85 remote portion of the process address space, an area
86 not usually available on a 10/50 system. Subsequent
87 10/50 type monitor calls cause a transfer to the com-
88 patability package which then interprets the call.

89 The compatibility routines are placed in the user
90 space for several reasons: (a) regular use can be made
91 of the pseudo interrupt system; (b) the compatability
92 package (which requires constant maintenance) can
93 be maintained as a separate module, totally independent
94 from the monitor; and (c) the monitor is protected
95 from malfunction by the compatibility routines.

96 4. User Interaction with TENEX 97 98 99

4.1 Terminal Interaction Capabilities 100

101 The terminal service module of TENEX was designed
102 to provide any type of interactive behavior a program
103 might find useful. Many programs, especially the com-
104 mand language interpreter described below, benefit
105 by having many short interactions with the user, often
106 one or a few characters. Full-duplex terminals are
107 preferred for use with TENEX for these reasons and for
108 the reason that the user can in fact anticipate the
109 machine's responses and begin typing input before
110 output is completed. Algorithms for echoing typein
111 ensure that the typescript is an accurate record of the
112 dialog. Half-duplex terminals may be used but at some
113 cost in convenience.

114 4.2. Executive Command Language 115

116 Users at terminals communicate and work with
117 TENEX primarily through a command language inter-
118 preter called the TENEX Executive, or EXEC. The EXEC
119

1 is an interactive, well human-engineered program
2 which can accept commands from a user's teletype
3 or from a file. It is implemented as a reentrant, shared
4 program which runs in user mode, usually as the top
5 level process in the structure.

6 The EXEC provides the user with a multitude of facil-
7 ities which are activated by simple, easy-to-learn com-
8 mands. These facilities provide access to the system
9 (e.g. LOGIN); utility operations on files and file direc-
10 tories; initiation of private programs and subsystems;
11 limited debugging aids; initiation of batch (detached
12 operations); printout of user information and system
13 statistics; and system maintenance.

14 The EXEC was designed with two primary objec-
15 tives—ease of learning and ease of use. To ease the
16 learning process, all commands are English words
17 which are descriptive of the facility being activated (e.g.
18 COPY to copy information from one file to another,
19 STATISTICS to obtain a listing of current system statis-
20 tics). Each command begins with a keyword. Depend-
21 ing on the command, the initial keyword may be
22 followed by arguments, such as file names, numbers,
23 and additional keywords, and/or “noise words” to
24 make the command more readable. The noise words
25 are enclosed in parentheses to distinguish them from
26 the arguments.

27 In order to help novice users, two special assistance
28 features were incorporated. First, when the EXEC
29 requires input from the user during a command inter-
30 action (for instance, to collect arguments of that
31 command), a cue is typed to indicate to the user what is
32 expected. For example, an interaction which renames
33 a file might be

34 @REN\$AME (EXISTING FILE) ALPHA\$.MAC
35 (TO BE) BETA

36 The user's input is bold face. The \$ indicates a typed ESC
37 (ASCII escape, code 33₈) which invokes the EXEC's
38 verbose cueing responses in parentheses. In this ex-
39 ample the user typed only three letters REN followed by
40 ESC which invoked command completion by the EXEC,
41 a feature which makes the language particularly easy
42 to use. An ESC after any initial substring of a command
43 or argument (such as a file name) invokes completion.
44 If the substring is insufficient for unique identification
45 of the intended input, the EXEC rings the teletype's
46 bell and awaits additional characters. If the initial
47 substring cannot be recognized the EXEC types “?”
48 to ask the user to retype that input. If the novice
49 user still doesn't understand what is expected in his
50 response, he may invoke the second special assistance
51 feature by typing the character “?”. This causes the
52 EXEC first to type out a list of all options available to
53 the user at that point and then request a response.
54
55
56
57
58
59
60

To summarize, three general styles of input may be
used, distinguished by syntactic structure, and including
special input terminators. Hence the styles do not
require different input modes, and thus may be inter-
mixed freely within a session or even within a statement,
adapting to the state of knowledge and verbosity of
the user. The input styles allow:

1. *Complete input.* A complete command may be
typed in, with all keywords and noise words given in
their entirety and without use of any nonprinting
characters.

2. *Abbreviations.* The user may shorten a command in
two ways: he can omit noise words completely and he
can shorten keywords. Any keyword may be abbrevi-
ated with any initial substring (terminated with space)
long enough to distinguish it from the other keywords
acceptable in that context.

3. *Completion.* The user types the same characters as
in abbreviated input, except he terminates each field
(keyword or argument) with the ESC key. This pro-
duces a print-out of the complete command—each
ESC causes the rest of the field (if an abbreviated key-
word or file name) and any following noise words
(with enclosing parentheses) to be printed.

The EXEC also provides editing characters to permit
the user to correct typing errors in his input. These
editing characters permit the user to delete the last
character of his typed input, the last word, or all of it.
He can also ask for his edited input to be retyped for
clarity.

4.3 Interrupt and Escape Characters

ASCII Control-C is the EXEC's attention character.
When typed by the user, it causes any running program
to be stopped and control to be given to the EXEC via
the pseudo interrupt system. The user may then con-
tinue his program or take any other action.

Another terminal interrupt character, Control-T,
is serviced by the EXEC. It interrupts a user's EXEC
process to type out the total CPU and console time
used and the status of the process being run under the
EXEC; this lower user process continues.

5. The Tenex File System

The TENEX file system provides a general mechanism
for obtaining information from and sending data to
external devices attached to the TENEX system [12].
Write-only and read-only devices are included in the
file system so that all TENEX I/O may be handled uni-
formly. The first major function of the TENEX file
system is to provide symbolic file name management.
This includes two separate but related activities. The
first involves translation of a symbolic name into an
internal “file descriptor block” pointer associated
with that name; the second involves checking informa-
tion concerned with (1) the file status, e.g. whether it

61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
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

1 exists, access rights, etc.; and (2) the process requesting
2 access to the file. This second activity, known as File
3 Access Protection, determines if this process should
4 be allowed to know about the existence of this file, and
5 if so, what access it is allowed.

6 A symbolic name for TENEX files consists of up to
7 five fields and thus conceptually represents a tree of
8 maximum depth five. Not all nodes of this tree go
9 down to maximum depth. This scheme was chosen
10 rather than a full tree to simplify the problem of com-
11 patibility with existing DEC PDP-10 software and name
12 lookup and recognition. We are currently considering
13 the feasibility of implementing a full tree directory
14 structure similar to MULTICS [11]. At each level there
15 would be a set of information, which is related to
16 access rights, and media dependence of the data
17 access for this node. Each node would represent a
18 collection of related information, with the terminal
19 nodes being files.

20 The fundamental unit of storage in a TENEX file is a
21 byte, which may be from 1 to 36 bits in length. A
22 stream of bytes constitutes a file, which is the basic
23 named element in the file system. Programs may refer-
24 ence files byte by byte in a sequential manner or, if the
25 device permits, at random. String (multiple byte)
26 transfers can also be made. No structure other than
27 bytes and files is imposed on the user, and byte and
28 string input and output are the basic operations. Of
29 course, additional structure and other operations may
30 be implemented by the user programs.

31 5.1 File Names

32 A TENEX file is named by a file descriptor composed
33 of five fields some of which are omitted for certain
34 devices. The five fields are device name, directory
35 name, file name, extension, and version number.

36 The file name field is intended to designate a class of
37 files which are related in some way. This convention is
38 not enforced, but most users of TENEX follow the con-
39 vention since it facilitates management of a user's
40 files. The extension field is intended to designate vari-
41 ously processed forms of the same information. A
42 file's extension is frequently specified by a program.
43 For example, PROG.MAC, PROG.REL, and PROG.SAV
44 would be used to indicate the MACRO assembly code
45 source, relocatable file, and binary image of a single
46 program.

47 The version number of a file enumerates successive
48 versions of a file. Normally each time a file is written a
49 new version is automatically created by making its
50 version number be one greater than the highest existing
51 version. This protects a user from loss if he acciden-
52 tally writes on the wrong file. Excess versions may be
53 deleted by the user, or automatically by the system,
54 when they have been put on a backup storage medium.

55 Any of the fields of a file description may be abbre-
56 viated except for device and version. The appearance
57 of an ESC in the file descriptor causes the portion of the

field before the ESC to be looked up, and the system
will supply the omitted characters and/or fields. Abbre-
viation without this output is not provided in order to
insure that the typescript reflects exactly what was
done. The system provides default values for each
field except the file name.

A default value is used for a field if the user omits
any input for that field, e.g. the device and directory.
This simplifies references to files in most common cases.

5.2 File Access Protection

Because TENEX must service a diverse user com-
munity, it is essential that access to files be protected
in a general way. Generally, access to a file depends on
two things: the kind of access desired and the relation
of the program making the access to the owner of the
file. Presently, a simple protection scheme is imple-
mented in which the only possible relationships a
program may bear to the file's owner are:

1. The directory attached to the job under which the
program is running is the same as the owning directory.
2. The directory attached to the job under which the
program is running is in the same group as the owning
directory.
3. Neither 1 nor 2.

Five kinds of access are distinguished for a file:
directory listing, read, write, execute, and append.
The above three relationships and five protection
types are related by 15 bits (a 3×5 binary matrix)
in which a one indicates that a particular access is
permitted for a particular relationship. If directory
listing access is not permitted, the process requesting
access is given an error return which is indistinguishable
from the error for nonexistent file. This is important
if the information that a file exists should not be gener-
ally available, as is the case for secure systems. Other
access restrictions cause errors only when an attempt
is made to open a file, as described below.

For purposes of determining group access, a 36 bit
word is administratively associated with each directory
and each user. If the bitwise "and" of the user group
word of the accessor and of the directory group word
of the accesser is nonzero, the group access permission
is used.

Provision has been made for a more general file
protection system in which more general access rela-
tionships may be expressed in a special file protection
language. For example, access may be allowed only
to an explicitly named set of users.

5.3 File Operations

Using a file in TENEX is basically a four step process:
first a correspondence is established between a file
name and a Job File Number (JFN), which is a small
index into a job table for files; next the file is opened,
establishing the mode and access permission and setting
up monitor tables to permit the data of the file to be
accessed; third, data is transferred to or from the file;

1 and finally, the file is closed, fixing up the directory
2 information and releasing the space occupied in system
3 tables for the file.

4 For purposes of file sharing, all instances of opening
5 a particular file should reference the same data. Data
6 written in a file will be immediately seen by readers
7 of the file. To protect against confusion resulting from
8 multiple uncooperating simultaneous writers and
9 readers of a file, a file can be opened with what we call
10 *thawed* or *unthawed* access. With thawed access, a
11 file may have any number of thawed writers and/or
12 thawed readers, but no provision is made to guarantee
13 that information is in a consistent (frozen) state.
14 With unthawed access, a file may have any number of
15 unthawed readers, or exactly one unthawed writer;
16 this prevents any potentially conflicting operations.
17 Simultaneous accessors of a file must be all thawed or
18 all unthawed.

20 6. The Monitor

22 6.1 Scheduler

23 The TENEX scheduler is designed to meet a set of
24 potentially conflicting requirements. The first is to
25 provide an equitable distribution of CPU service, which
26 we define as at least $1/N$ of real time where there are N
27 jobs on the system. Secondly, because TENEX is de-
28 signed to be a good interactive system, the scheduler
29 must identify and give prompt service to jobs making
30 interactive requests. Thirdly, because use of the CPU
31 is intimately tied to the allocation of core memory,
32 it must make efficient use of core memory to maximize
33 CPU usage. Finally, the scheduler should have provision
34 for administratively controlling the allocation of re-
35 sources so as to obtain other than equal distribution if
36 desired.

37 **6.1.1 Balance Set Scheduling.** For the scheduler, we
38 want a coherent policy which obeys Denning's "work-
39 ing set principle" [6]. A priority rating, as described
40 below, is given to each runnable process in the system,
41 and an estimate is made of the working set size of each
42 process. The jobs with highest priority whose total
43 working sets will fit in core are called the balance set
44 and may be run concurrently. When any process in the
45 set page faults one of the others in the balance set is
46 given CPU service. Denning [7, 8] has shown that such
47 balance set scheduling minimizes thrashing and tends
48 to maximize system efficiency. Periodic monitoring of
49 the entire set of runnable processes for changes in
50 priorities and in working set sizes allows adjustment
51 of balance set membership.

52 **6.1.2 Setting Process Priorities.** To implement the
53 basic scheduling function, a scheduling algorithm was
54 chosen which groups processes together on a number of
55 separate queues each with an associated runtime quan-
56 tum, similar to algorithms described by Corbato [3]
57 and BBN [1]. Lower queues in general have lower prior-

ities but longer runtimes. A common problem with
many schedulers of this type is that processes are
placed on the highest priority queue after any interac-
tion. Under conditions of heavy load or with poorly
behaved interactive processes, it may happen that the
interactive processes succeed in using all of the avail-
able time and so lock out the compute-bound processes
which have fallen to the lower queues.

In TENEX, priority is based on a long term average
ratio of CPU use to real time, and a process's priority
after an interaction is determined by its priority before
the interaction and the length of the interaction.
Specifically, a process's priority is decreased while
running at a constant rate, C , and increased while
blocked at a rate of C/N , where N is the number of
runnable processes in the system. This ensures that
equitable service is given both to compute-bound and
interactive jobs.

To improve response characteristics, an interactive
"escape clause" is included in the scheduling algorithm.
After a block wait of greater than minimum time, a
process is given a short quantum at maximum priority.
Priority and queue position after this burst are deter-
mined by the long term average. The effect of this
provision is to ensure quick service to very short inter-
actions, even when requested immediately after a
long computation.

6.1.3 Resource Guarantees and Limitations. In
some cases CPU resource guarantees independent of
load are desired, e.g. a demonstration which requires
significant CPU time during a period of medium or
heavy load, or a user who is willing to pay extra for
premium service which does not degrade as the load
on the machine increases. A facility is implemented in
TENEX to handle these situations.

A person with appropriate administrative access
can assign to any job on the system a fraction, F , of
guaranteed CPU service. For any job so designated, the
scheduler will attempt to ensure that

$$C/T > F,$$

where C is the CPU seconds used by the process, and
 T is the real time since the process last unblocked.
For example, if the parameter is set to 30 percent,
the scheduler will provide at least 18 seconds of CPU
service to the specified job during each minute of real
time.

This parameter acts as a ceiling as well as a floor
for CPU service. That is, if there are other runnable
processes on the system which are not declared special,
then the scheduler will ensure that the special process
receives no *more* than the stated fraction of CPU service.
We have found that a process with this sort of re-
source guarantee displays very consistent interactive
behavior despite widely varying loads on the time
sharing machine.

61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
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

6.2 Core Management

The information provided in the core status table by the paging hardware is essential to the proper management of core memory in TENEX to avoid thrashing and other forms of inefficient operation. Paging is done on demand. No ordinary pages are preloaded before a process is run, and in general, a process will not have all the pages of its virtual memory in core at once.

When a process references a page which is not in core, a pager trap occurs and a core management routine is invoked. The run time since the last page fault is used for a running average of page fault times, i.e. interfault intervals in process time. A process is considered to have enough of its working set if its average page fault time equals PAV, a system parameter currently set to 67ms (or 2 drum revolutions). If the process is faulting more often than PAV, it is considered below its working set size, and a swap to bring in the requested page initiated. Control returns to the scheduler so that it can run another process until the swap is complete. If the process is faulting less often than PAV, a core management routine is invoked to reduce the size of the process, i.e. remove some of its pages from core if space is needed. We have found this algorithm is stable and works well; Denning and Schwartz [9] show formally that this is a good idea.

To reduce the size of a process working set, a least recently used algorithm is used. The age of each page of a process is determined by a 9 bit logical age field stored by the pager in the core status table when a reference to that page causes a pager reload. Since collecting pages is costly, all sufficiently old pages are removed on a working set reduction. A quick reference however can catch a page before it leaves physical core.

6.3 System Measurements

In order to observe and improve the performance of TENEX in regular service, various measuring functions were built into monitor routines. Some measures serve to indicate the efficiency of scheduling, the core/CPU balance, and the nature of the various processes running on the system. The scheduler maintains a set of integrals over time which give (as a fraction of real time):

IDLE, time when no processes are requesting CPU service.

WAIT, time when all runnable processes are waiting for completion of page fault.

CORE, overhead time spent in core management.

TRAP, time spent handling pager traps.

Two relationships among these are:

$s = \text{IDLE} + \text{WAIT} + \text{CORE} = \text{total time in scheduler.}$

$\text{REALTIME} - s - \text{TRAP} = \text{time spent running user processes.}$

Also maintained by the scheduler is an integral over time of the number of processes in the balance set, the number of transfers between core and second-

ary storage, and the number of terminal interactions. 61

One measure is of interest on a recurring basis to all users of the system. The scheduler maintains three exponential averages 62 63 64

$$A(T+t) = A(T) * \exp(-t/c) + N * (1 - \exp(-t/c)),$$

with time constants c equal to 1, 5, 15 minutes and N equal to the number of runnable processes on the system. This indicates the true current load on the system better than the number of jobs logged in. Users often choose on the basis of these load figures what they do on the system at a particular time. Three figures are better than just the last one, because from these the user can predict the trend as well as note local variation. 65 66 67 68 69 70 71 72 73 74 75 76

6.4 Debugging Aids

Certain debugging procedures and aids used in the development of the system contributed greatly to the speed of development and integrity of the system. Our principal debugging aid is a program called DDT, available in several forms in the system. DDT is a program which allows memory locations to be examined and modified, and breakpoints (return of control to DDT) to be placed in a running program. All interactions with DDT are symbolic, using the symbols defined in the source program and obtained from the assembler. 77 78 79 80 81 82 83 84 85 86 87 88

The form of DDT first used and still necessary for debugging basic level code is a stand-alone version which resides in core memory along with monitor. It is used for debugging the scheduler, portions of the core manager, and other basic routines. The second form of DDT was added as soon as the basic monitor could support demand paging and create a virtual memory. This DDT exists in the monitor map and may be used as an ordinary program at a system terminal. It is capable of examining and changing the running monitor and all of the associated tables and other contents of the monitor virtual memory. Use of this form of DDT actually allows several persons to work on debugging portions of the monitor simultaneously. A third form of DDT is used with user programs and is cognizant of the access status (execute or write protection, etc.) of pages of the user program. 89 90 91 92 93 94 95 96 97 98 99 100 101 102 103 104 105

Debugging the system was further facilitated by the use of a considerable number of internal redundancy and consistency checks which call one of two recovery routines upon detecting any malfunction. If the system is attended by system personnel, these routines enter a DDT breakpoint and the state of the monitor can be examined to determine what has gone wrong. This has enabled us to find and correct the obscure or infrequent faults in the software. If the system is unattended, then, depending on which routine was called and the severity of the fault, one of three things happens: (1) operation is continued with no interruption; (2) one process is crashed; or (3) the 106 107 108 109 110 111 112 113 114 115 116 117 118 119

1 system is automatically reloaded and restarted. This
2 usually provides continued operation in the event of
3 unexpected hardware or software malfunction. Both
4 routines cause the source of the inconsistency and a
5 message describing the problem to be logged for fur-
6 ther action by system personnel.

7 7. Conclusion

10 One of the most valuable results of our work was
11 the knowledge we gained of how to organize a hard-
12 ware/software project of this size. Virtually all of the
13 work on TENEX from initial inception to a useable sys-
14 tem was done over a two year period. There were a
15 total of six people principally involved in the design
16 and implementation. An 18 month part-time study,
17 hardware design and implementation culminated in a
18 series of documents which describe in considerable
19 detail each of the important modules of the system.
20 These documents were carefully and closely followed
21 during the actual coding of the system. The first stage
22 of coding was completed in 6 months; at this point the
23 system was operating and capable of sustaining use by
24 nonsystem users for work on their individual projects.
25 The key design document, the JSYS Manual (extended
26 machine code), was kept updated by a person who
27 devoted full time to insuring its consistency and co-
28 herence; and in retrospect, it is our judgment that this
29 contributed significantly to the overall integrity of the
30 system.

32 We felt it was extremely important to optimize the
33 size of the tasks and the number of people working
34 on the project. We felt that too many people working
35 on a particular task or too great an overlap of people
36 on separate tasks would result in serious inefficiency.
37 Therefore, tasks given to each person were as large as
38 could reasonably be handled by that person, and insofar
39 as possible, tasks were independent or related in ways
40 that were well defined and documented. We believe that
41 this procedure was a major factor in the demonstrated
42 integrity of the system as well as in the speed with which
43 it was implemented.

45 *Acknowledgments.* In addition to the work done by
46 the authors, significant contributions to the design
47 and implementation of TENEX were made by T.R.
48 Strölo, who led the technical staff, and J.R. Barnaby,
49 who implemented the EXEC subsystem. Others who
50 contributed are T. Myer, E. Fiala, D. Wallace, and
51 J. Elkind.

References

1. BBN Medical Information Technology Department. The hospital computer project time sharing executive system, BBN Rep. No. 1673, Apr. 1968.
2. Bobrow, D.G., Burchfiel, J.D., Murphy, D.L., and Tomlinson, R.S. TENEX, a paged time sharing system for the PDP-10. BBN Rep. No. 2180, Aug. 1971.
3. Corbató, F.J., et al. An experimental time-sharing system. Proc. AFIPS 1962 SJCC, Vol. 21 Spartan Books, New York, pp. 335-344.
4. Corbató, F.J., et al. An introduction and overview of the Multics system. Proc. AFIPS 1965 FJCC, Vol. 27 Pt. 1, Spartan Books, New York, pp. 185-196.
5. Digital Equipment Corp. PDP-10 Reference Handbook. Dec. 1971.
6. Denning, P. The working set model for program behavior. *Comm. ACM* 11, 5 (May 1968), 323-333.
7. Denning, P. Thrashing, it's causes and prevention. Proc. AFIPS 1968 FJCC, Vol. 33 Pt. 1, AFIPS Press, Montvale, N.J., pp. 915-922.
8. Denning, P. Equipment configuration in balanced computer systems. *IEEE Trans. Comput. C-18*, 11 (Nov. 1969), 1008-1012.
9. Denning, P., and Schwartz, S.C. Properties of the working set model. *Comm. ACM* 15, 3 (Mar. 1972), 189-196.
10. Lampson, B., et al. A user machine in a time sharing system. *Proc. IEEE* 54, 12, (Dec. 1966), 1766-1774.
11. Spier, M.J., and Organick, E. The Multics interprocess communication facility. Proc. Sec. Symp. Oper. Sys. Princ., Oct. 1969, ACM, New York, pp. 83-91.
12. Wilkes, M. *Time Sharing Computer Systems*. American Elsevier, New York, 1968.

61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
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