Title- A Structured Operating System
for a PDP-11/45

Other Keywords- UNIX
Multi-Environment
Real-Time

| Author(s) | Location and Room | Extension | Charging Case- 39394 |
|---|---|---|---|
| Lycklama, H. | MH 7C-211 | 6170 | |
| Bayer, D. L. | MH 7C-207 | 3080 | Filing Case- 39394-11 |

## ABSTRACT

A structured operating system, MERT, consisting of
a set of autonomous processes has been designed and implemented
on a PDP-11/45 computer in Department 1352. The MERT system is
a multi-environment, real-time operating system consisting of a
set of basic kernel procedures providing services for the inde-
pendent processes which support the different operating system
environments. A well-developed set of inter-process communica-
tion primitives have been implemented, including event flags,
message buffers, shared memory and shared files. We believe
this provides a good base for providing support for various
operating system environments and for providing real-time
response for processes. The UNIX time-sharing system has been
implemented as one environment on the MERT system. This paper
(*) provides an overview of the MERT system.

* - Submitted to COMPCON Conference, September 7-9, 1975

DISTRIBUTION
(REFER GEI 13.9-3)

| OMPLETE MEMORANDUM TO | COMPLETE MEMORANDUM TO | COVER SHEET ONLY TO | COVER SHEET ONLY TO | COVER SHEET ONLY TO |
|---|---|---|---|---|
| ORRESPONDENCE FILES | LOZIER,JOHN C | COVER SHEET ONLY TO | CHANG,HERBERT Y | FULTON,ALAN W |
| | LUDERER,GOTTFRIED W R | | CHANG,S-J | GARCIA,R F |
| OFFICIAL FILE COPY | LYCKLAMA,HEINZ | | CHAPPELL,S G | GAY,FRANCIS A |
| PLUS ONE COPY FOR | LYONS,T G | CORRESPONDENCE FILES | CHEN,STEPHEN | GEER,EUGENE W JR |
| EACH ADDITIONAL FILING | MACHOL,R E JR | | CHEN,T L | GELLINEAU,A C |
| CASE REFERENCED | MAKER,MS D A | 4 COPIES PLUS ONE | CHERRY,MS L L | GELLIS,H S |
| | MALTHANER,W A | COPY FOR EACH FILING | CHIANG,T C | GEPNER,JAMES R |
| DATE FILE COPY | MARANZANO,JOSEPH F | CASE | CHODROW,MARK M | GEYLING,F T |
| (FORM E-1328) | MASHEY,JOHN R | | CHRIST,C W JR | GIBB,KENNETH R |
| | MC ILROY,M DOUGLAS | | CIRILLO,CARL | GILBERT,MRS HINDA S |
| 10 REFERENCE COPIES | MCDONALD,H S | ABRAHAM,STUART A | CLAYTON,D P | GIMPEL,JAMES F |
| | MILLER,S E | ACKERMAN,A F | CLIFFORD,ROBERT M | GITHENS,JOHN A |
| LBERTS,BARBARA A | MORGAN,S P | AHO,A V | CLOUTIER,J E | GITLIN,RICHARD D |
| NDERSON,MRS C M | NINKE,WILLIAM H | AHRENS,RAINER B | COBEN,ROBERT M | GLUCK,F |
| RDIS,R B | O CONNELL,T F | ALCALAY,DAVID | COHEN,HARVEY | GOETZ,FRANK M |
| RNOLD,S L | O NEILL,DENNIS M | ALLEN,JAMES R | COLDREN,LARRY A | GOGUEN,MS NANCY |
| AYER,DOUGLAS L | OSSANNA,J F JR | ALLES,HAROLD G | COLE,LOUIS M | GOLABEK,MISS R |
| IREN,MRS IRMA B | PATEL,C K N | ALMQUIST,R P | COLE,M O | GOLDSTEIN,A JAY |
| LUM,MRS MARION | PERDUE,R J | AMORY,R W | COLLIER,ROBERT J | GORDON,P L |
| OYD,GARY D | PEREZ,MRS CATHERINE D | AMOSS,JOHN J | COLTON,JOHN R | GRAHAM,R L |
| RANDT,RICHARD B | PERNESKI,A J | ANDERSON,M M | COPP,DAVID H | GRAMPP,F T |
| ROWN,W STANLEY | PETERSON,RALPH W | ARNOLD,GEORGE W | COSTANTINO,B B | GREENBAUM,H J |
| UCHSBAUM,S J | PHILLIPS,S J | ARTHURS,EDWARD | COSTON,WALTER P | GREENE,MRS DELTA A |
| URROWS,T A | PILLA,MICHAEL A | ATAL,B S | COULTER,J REGINALD | GREENHALGH,H WAIN |
| ANADAY,RUDD H | PINSON,ELLIOT N | BAKER,BRENDA A | COURTNEY PRATT,J S | GROSS,ARTHUR G |
| ARDOZA,WAYNE M | PLAUGER,P J | BALDWIN,GARY L | CRUME,LARRY L | GUERRIERO,JOSEPH R |
| ARRAN,J H | POPPER,C | BALDWIN,GEORGE L | D STEFAN,D J | HAFER,E H |
| ARR,DAVID C | +PRIM,ROBERT C | BARTLETT,WADE S | DAVIDSON,CHARLES L | HAGELBARGER,D W |
| HRISTENSEN,C | ROBERTS,CHARLES S | BASEIL,RICHARD J | DETRANO,MRS M K | HAGGERTY,J F |
| LOGSTON,A M | ROCHKIND,M J | BAUER,BARBARA T | DEUTSCH,DAVID N | HAHN,J R JR |
| ONDON,J H | RODIHAN,MRS PATRICIA A | BAUGH,C R | DICKMAN,B N | HALFIN,SHLOMO |
| OOK,THOMAS J | ROSLER,LAWRENCE | BECKER,R A | DIMMICK,JAMES O | HALL,ANDREW D JR |
| RANE,RODERICK P | SATZ,L R | BECKETT,J T | DOMPIERRE,J A | HALL,MILTON S JR |
| UNNINGHAM,STEPHEN J | SIX,FREDERICK B | BERGLAND,G DAVID | DONOFRIO,L J | HAMILTON,PATRICIA |
| UTLER,C CHAPIN | SLICHTER,W P | BERNSTEIN,LAWRENCE | DREIZLER,HOWARD K | HANSEN,MRS G J |
| E JAGER,D S | SMITH,D W | BEYER,JEAN-DAVID | DRISCOLL,PATRICK J | HARRISON,NEAL T |
| ICK,GEORGE W | SPENCE,NORMAN A | BILOWOS,RICHARD M | EDELSON,D | HARTWELL,WALTER T |
| OLOTTA,T A | STAMPFEL,JOHN P | BIRCHALL,R H | EITELBACH,DAVID L | HARUTA,K |
| OWD,PATRICK G | STEVENSON,D E | BLEICHER,EDWIN | ELLIOTT,R J | HASZTO,EDWARD D |
| DMUNDS,T W | STURMAN,JOEL N | BLINN,JAMES C | ELY,T C | HAUSE,A D |
| RRICHIELLO,PHILIP M | SWANSON,GEORGE K | BLY,JOSEPH A | ESSERMAN,ALAN R | HEATH,SIDNEY F III |
| FEDER,J | TAGUE,BERKLEY A | BODEN,F J | ESTOCK,R G | HELD,RICHARD W |
| FORTNEY,MRS VIRGINIA J | TERRY,M E | BOHACHEVSKY,I O | FABISCH,MICHAEL P | HENIG,MRS FRANCES H |
| FRANK,H G | TEWKSBURY,S K | BOURNE,STEPHEN R | FARGO,GEORGE A | HERGENHAN,C B |
| FREENY,S L | THOMPSON,JOHN S | BOWERS,J L | FELS,ALLEN M | HERMAN,KENNETH M |
| ATES,G W | THURSTON,R N | BOWYER,L RAY | FIGLIUZZI,MISS M E | HEROLD,JOHN W |
| ILLETTE,DEAN | TILLOTSON,L C | BOYCE,W M | FIORE,MRS RHODA J | HESS,MILTON S |
| IORDANO,PHILIP P | UNDERWOOD,R W | BRAINARD,RALPH C | FISCHER,H B | HOLTMAN,JAMES P |
| LASSER,ALAN L | VIGGIANO,F A | BREECE,HARRY T III | FLANAGAN,J L | HONIG,W L |
| RAVEMAN,R F | VOGEL,G C | BREITHAUPT,ALLAN R | FLEISCHER,HERBERT I | HOPPERT,D J |
| HAIGHT,R C | WALKER,MISS E A | BROWN,COLIN W | FLUHR,ZACHARY C | HORNBACH,THOMAS S |
| HAMMING,R W | WANDZILAK,P D | BUTLETT,D L | FOUGHT,B T | HOYT,WILLIAM F |
| HANNAY,N B | WATKINS,G T | BUTZIEN,PAUL E | FOUNTOUKIDIS,A | HUDSON,E T |
| HASKELL,BARRY G | WEBB,FRANCIS J | BYRNE,EDWARD R | FOWLER,BRUCE R | HUMCKE,D J |
| HUPKA,MRS FLORENCE | WEHR,L A | BZOWY,D E | FOY,J C | HUNNICUTT,CHARLES F |
| IVIE,EVAN L | WELLER,DAVID R | CABLE,GORDON G JR | FRANK,MISS A J | IMAGNA,CLYDE P |
| JOHNSON,STEPHEN C | WHITE,RALPH C JR | CAMPBELL,J H | FRANK,RUDOLPH J | IPPOLITI,O D |
| KAPLAN,A E | WILSON,GEOFFREY A | CANDY,JAMES C | FRASER,A G | IRVINE,M M |
| KAUFMAN,LARRY S | WOOD,J L | CASEY,JOSEPH P | FREEDMAN,M I | JACKOWSKI,D J |
| KEVORKIAN,DOUGLAS E | YAMIN,MRS E E | CASPERS,MRS BARBARA E | FREEMAN,K GLENN | JACOBS,H S |
| LARSEN,ARTHUR B | YOUNG,JAMES A | CAVINESS,JOHN D | FREEMAN,R DON | JAMES,DENNIS B |
| LESSEK,PETER V | 112 NAMES | CHAMBERS,J M | FREIDENREICH,MRS B | JARVIS,JOHN F |
| LIMB,J O | | CHAMBERS,MRS B C | FROST,H BONNELL | JENKINS,J MICHAEL |

+ NAMED BY AUTHOR                    > CITED AS REFERENCE SOURCE

...
438 TOTAL

ERCURY SPECIFICATION...................................................................................................

OMPLETE MEMO TO:
  135-DPH        13-DIR        11-EXD        15-EXD        16-EXD        127-SUP        135-SUP
UNOS   = UNIX/OPERATING SYSTEM

OVER SHEET TO:
  135-MTS       9152-MTS    1271   1273   8234
COOS   = COMPUTING/OPERATING SYSTEMS/SURVEY PAPERS ONLY

------------------------------------------------------------------------------------------------------------

'O GET A COMPLETE COPY:

. BE SURE YOUR CORRECT ADDRESS IS GIVEN ON THE OTHER SIDE.
:. FOLD THIS SHEET IN HALF WITH THIS SIDE OUT AND STAPLE.
. CIRCLE THE ADDRESS AT RIGHT.   USE NO ENVELOPE.

PLEASE SEND A COMPLETE COPY TO THE ADDRESS SHOWN ON THE
OTHER SIDE
NO ENVELOPE WILL BE NEEDED IF YOU SIMPLY STAPLE THIS COVER
SHEET TO THE COMPLETE COPY.
IF COPIES ARE NO LONGER AVAILABLE PLEASE FORWARD THIS
REQUEST TO THE CORRESPONDENCE FILES.

Memorandum for File

## Introduction

As operating systems become more sophisticated and complex, providing more and more services for the user, they become increasingly difficult to modify and maintain. Fixing a "bug" in some part of the system may very likely introduce another "bug" in a seemingly unrelated section of code. Changing a data structure is likely to have major impact on the total system. It has thus become increasingly apparent over the past years that adhering to the principals of structured modularity (1),(2) is the correct approach to building an operating system. The influence of a process must be confined to an environment which is well protected from the rest of the system and must never affect the state of other environments.

Brinch Hansen (2) implemented a nucleus of a multi-programming system on an RC 4000 computer using message buffering

as the basic means of inter-process communication. Our system uses a different set of message buffering primitives as well as other communication primitives to achieve process synchronization and information transfer in an efficient manner.

The Hardware

The PDP-11/45 computer provides an eight-level hierarchical interrupt structure with priority levels numbered from 0 (lowest) to 7 (highest). Associated with the interrupt structure is the programmed interrupt register which permits the processor to generate interrupts at priorities of one through seven. The programmed interrupt serves as the basic mechanism for driving the system.

The memory management unit provides a separate set of address mapping and access control registers for each of the processor modes: kernel, supervisor and user. Furthermore, each virtual address space can provide separate maps for instruction references (called I-space) and data references (called D-space). Segments may be given read/write protection (3). Each process has access only to its own separate address space.

System Structure

The basic computer hardware resources consist of the actual memory, the CPU and the various I/O devices. The first level (see Fig. 1) of the operating system structure, called the kernel, controls and allocates these resources. The kernel consists of a set of highly privileged procedures and therefore must be

very reliable.

The second level of software consists of kernel-mode processes which comprise the various I/O device drivers. Each process at this level has access to a limited number of I-space base registers in the kernel mode, providing a firewall between it and sensitive system data accessible only using D-space mode.

At the third software level are the various operating system supervisors which run in supervisor mode. These processes provide the environments which the user sees and the interface to the basic kernel services.

At the fourth level are the various user procedures which execute in user mode under control of the supervisory environments. The primitives available to the user are provided by the supervisory environments which catch the user traps. Actually the user procedure is merely an extension of the supervisor process. This is the highest level of protection provided by the PDP-11/45 hardware. Higher levels of software such as application packages can achieve protection through software strategies such as running in interpretive mode.

## Processes and Segments

A process is a related collection of logical segments which are executed by the processor. A logical segment is a piece of contiguous memory, 32 to 32K 16-bit words, which can grow in increments of 32 words. A name may be associated with a segment

to provide a label for sharing it between independent processes. This name is derived from the file system and is thus globally known to all processes.

The capabilities of a process are defined by the mode of the processor while executing the process segments. A kernel process, running in kernel mode, is limited to 12K words, is locked in memory and has direct access to the I/O registers. Kernel processes can execute at processor priorities of three to seven and are driven by both software and external interrupts. These processes can respond to interrupts within 100 microseconds.

All processes which do not run in kernel mode are supervisor processes. These processes all execute at processor priorities one and zero and are scheduled by the kernel scheduler process. Associated with each supervisor process is a special read-only segment (in supervisor mode) called the process control block (PCB). The PCB contains space to save the state of the process and a segment table which describes the virtual address space of the process. The process can manipulate the segment table by trapping to the kernel. Supervisor processes can be either swappable or non-swap depending upon real time constraints. The maximum size of a process is limited by the addressing capability of the hardware to 124K words. Real time response attainable by supervisor processes is of the order of several milliseconds.

## Process Communication

A number of methods are available to the processes for communicating with each other. These include:

```
event flags
messages
shared memory
files
```

These methods can be used to guarantee the synchronization of parallel, cooperating processes and for maintaining the modularity of these processes.

The most basic communication mechanism transmits one bit event flags:

```
event(procid, event)
```

which will cause the process procid to receive an interrupt and the word event, which contains one or more event flags. The receiving process can inhibit event interrupts by selectively enabling particular events.

A higher-bandwidth communication path is provided by the message sending and receiving mechanisms in the kernel. A common pool of message buffers (multiple of 16 words) and a message queue is provided for each process. See Figure 2 for the message format. Communication between all processes can be achieved by means of the following kernel primitives:

```
send message(message-buffer)

receive message(message-buffer)

receive message type(message-buffer)
```

Send message copies the message from the process message buffer into the next available kernel message buffer and places it on the input message queue of the intended receiver and sends a message event to the receiver. Receive message copies a message

into the process message buffer if there is a message on its
input queue. If the queue is empty, the process decides whether
to continue or roadblock while waiting for a message. A process
may also look for a particular type of message on its input mes-
sage queue by means of the receive message type primitive. This
enables a process to synchronize the processing of its ack-
nowledgement messages. The transmission of a message almost
always results in the receipt of an acknowledgement or answer
message.

Processes have the capability to share segments. For exam-
ple, a segment may be shared as a common pool of I/O buffers for
the use of co-operating processes. This sharing mechanism is
used in the implementation of a time-sharing supervisor in the
system. Shared segments are also be used to transfer information
between two co-operating processes.

The use of files for inter-process communication is common
in many operating systems. In our case named files are passed
between communicating processes. We also provide the mechanism
to lock a file, i.e. make it inaccessible to other processes.
The existence of the locked file is a means of inter-process com-
munication. The contents of a file may of course also be shared
by a number of processes.

## The Kernel

The kernel consists of a process dispatcher, a trap handler,
and routines which implement the system primitives.

The process dispatcher fields all interrupts, sets up the appropriate kernel mode process, and transfers control to it.

The trap handler fields all traps and faults and, in most cases, transfers control to a trap handling routine in the process which caused the trap or fault. For the purposes of debugging, the "break point trap" executed from supervisor or kernel mode will cause an image of the process to be written in a file and the process to be terminated. This feature has proven to be an invaluable tool for debugging kernel processes which control I/O devices and new versions of our time sharing supervisor.

The kernel primitives deal primarily with four logical functions:

1) The sending and receiving of messages and events.

2) The allocation and manipulation of segments which includes such operations as locking for I/O, growing, and sharing segments.

3) The manipulation of a process virtual address space. This includes specification of access permissions, starting virtual address, and growth direction of segments.

4) The manipulation of process attributes which influence scheduling. Included here would be altering the process priority and making the segments of a process unswappable.

Closely associated with the kernel are the memory management and scheduler processes. These two processes are special in the sense that they have access to kernel D-space and reside in the kernel segments. In all other respects they follow the

discipline established for kernel processes.

The memory manager process communicates with the rest of the system via messages and is capable of handling three types of requests:

1) Setting the segments of a process into the active state, making space by swapping or shifting other segments if necessary.

2) Loading and locking a segment contiguous with other locked segments to reduce memory fragmentation.

3) Deactivating the segments of a process.

The scheduler process is responsible for scheduling all supervisor mode processes. The scheduler utilizes both time-sliced, round robin and preemptive priority scheduling techniques. The main responsibility of the scheduler is to select the next process to be executed. The actual loading of the process is accomplished by the memory manager.

## File System

The system has the capability to provide such diverse environments as those suitable for real-time applications, data-base query systems and time-sharing systems. The types of files normally provided by operating systems cover the range from linked files to large and contiguous files, although these types rarely all appear on one system. Large and contiguous files in particular are appropriate for real-time applications and data-base query systems, and lend themselves to more sophisticated file I/O capabilities.

The file manager process provides the capabilities for all of these file types using a single type of file map entry. A file map entry describes the file layout by a series of 2-word extents consisting of starting block number and number of consecutive blocks. This format enables the user to allocate secondary storage space for a large contiguous file. For time-sharing applications the file may consist of a number of small extents. The maximum file size which can currently be accommodated under this scheme is 850,000,000 bytes.

The file manager process is implemented as a separate process controlling all accesses to files and is the only process which may alter the file map entries. Other processes communicate with the file manager entirely by means of messages. The messages to the file manager are handled by a number of tasks running in parallel, one on behalf of each message.

Status

The Operating System as discussed above has been implemented on a PDP-11/45 computer with 80K words of memory. The size of the kernel is 8K including system tables. Two basic system processes are started up with the system, a process manager and a file manager. A time-sharing supervisor, logically equivalent to UNIX (4), has been implemented as an environment, requiring about 11K of memory. At the current state of implementation of the system, the system overhead introduced by the message traffic and context changing reduces the throughput of the UNIX time-sharing

processes by about 5 percent.

Development of new and existing supervisors is being done using this time-sharing system for which many user application programs have been written over the years. One real-time supervisor currently being developed is an experimental telephone office common control used for call processing. This process is locked in memory and runs at high priority upon receiving a message or an event. This process may be loaded dynamically and run under control of the time-sharing supervisor. It is debugged by planting break-point traps and dumping out a memory image of the process on a file. In fact many of the kernel device drivers have been debugged in this manner. This allows the system to remain operational while some very sensitive code is being debugged. The debugging facilities allow a user to examine the state of any process in detail either while running or after being dumped.

The flexibility of the system is perhaps best shown by the ability to run and debug a new version of the time-sharing supervisor using the current version of that supervisor. This enables the system designer to test out new features of the time-sharing supervisor without running the risk of disrupting other users. The debugged version of the supervisor may then be installed without taking down the system. The structure of the operating system continues to provide a base for doing further operating systems research.
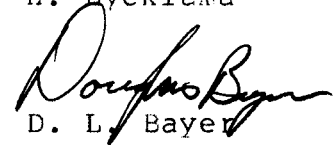
## Acknowledgments

Some of the concepts incorporated in the basic kernel were developed in a previous design and implementation of an operating system kernel by Mr. C. S. Roberts and one of the authors (H. Lycklama). The authors are pleased to acknowledge Mr. C. S. Roberts for many fruitful discussions during the design stage of the current operating system.

H. Lycklama

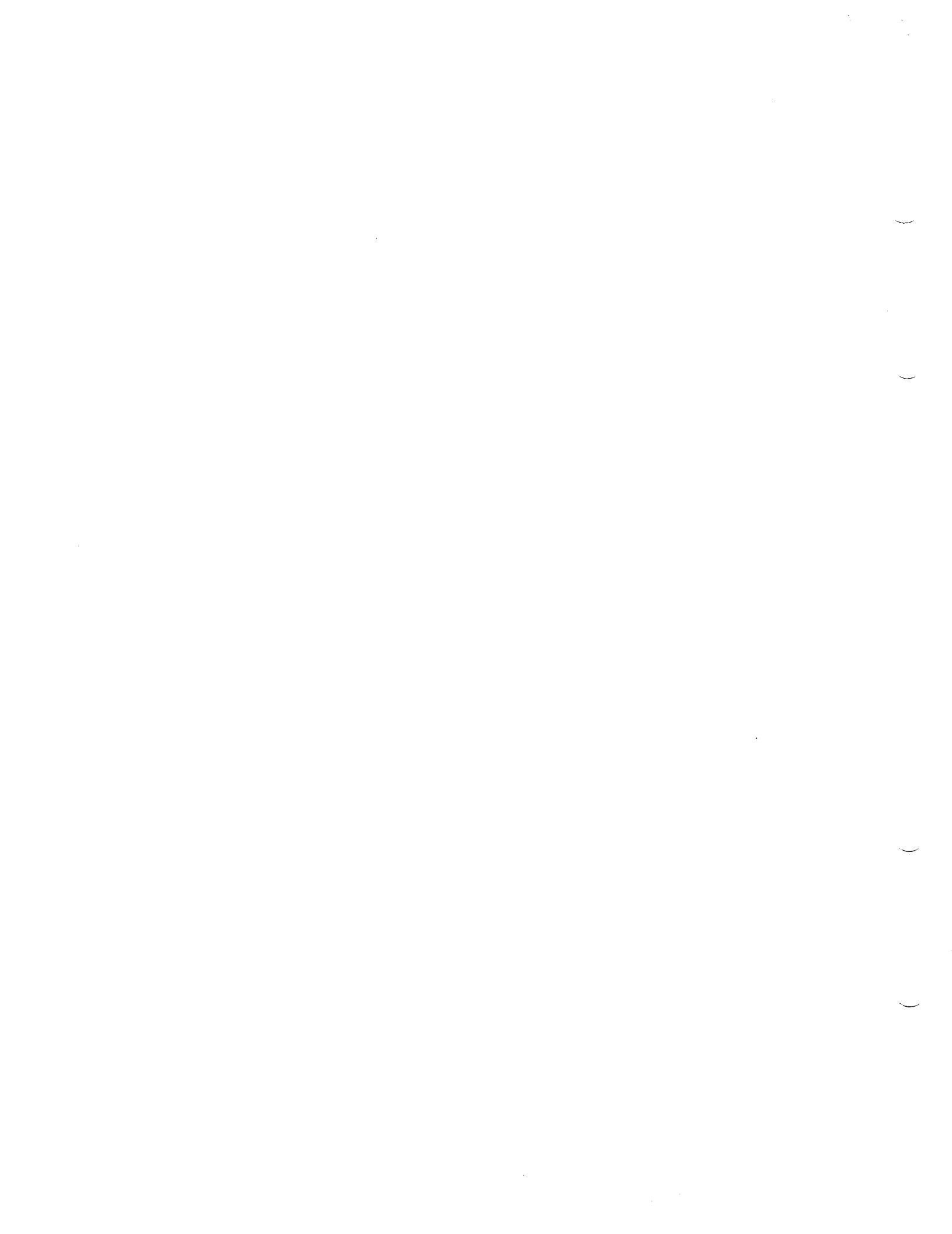D. L. Bayer

MH-1352-HL-JER

Att.
References
Figures

## References

(1) Dijkstra, E.W., The Structure of the 'THE' Multi-Programming System.  Comm. ACM 11, (May 1968), p341.

(2) Brincn Hansen, P., The Nucleus of a Multi-Programming System. Comm. ACM 13, (April 1970), p238.

(3) PDP -11/45 Processor Handbook, Digital Equipment Corporation, Maynard MA, 1971.

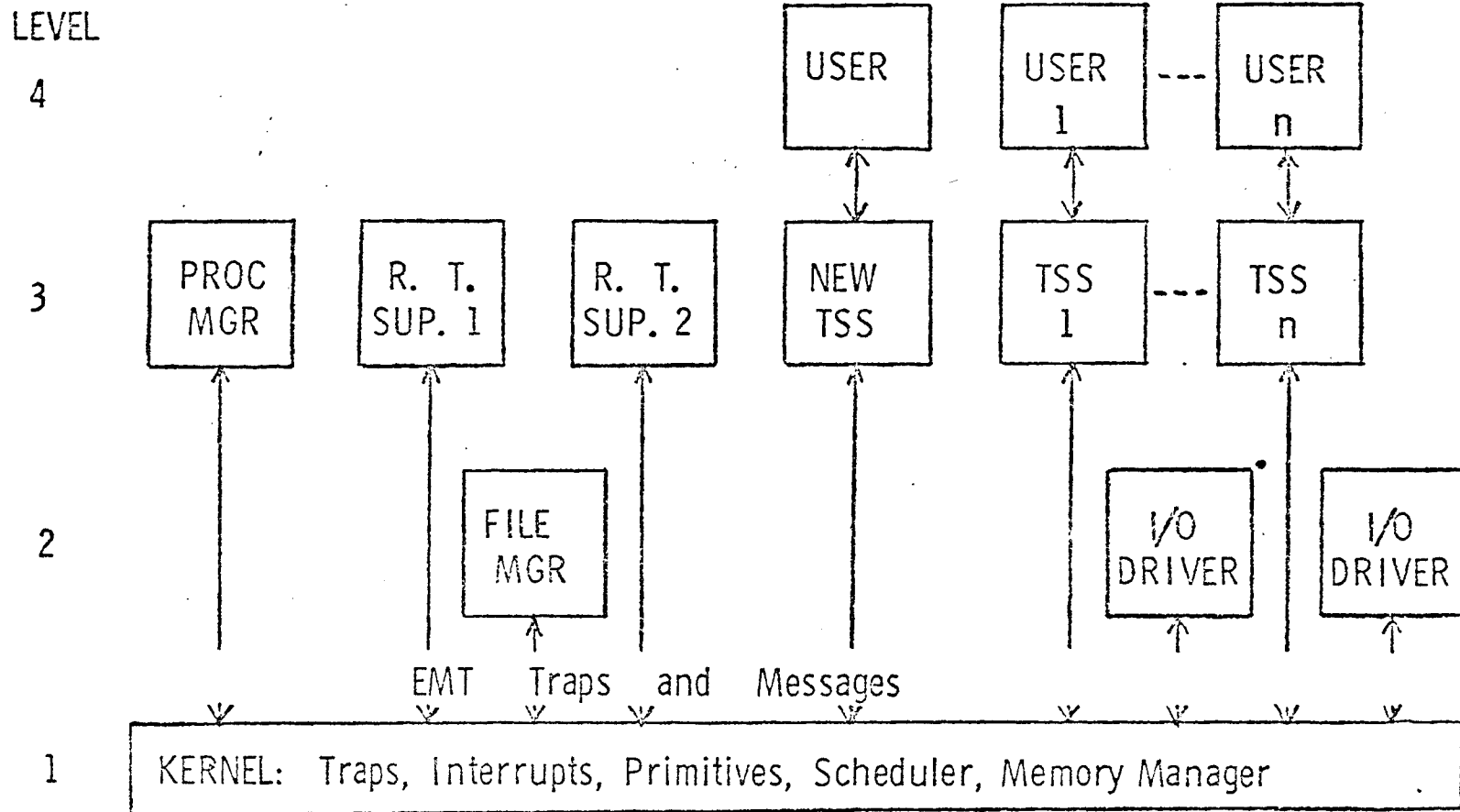(4) Thompson, K. and Ritchie, D.M., The UNIX Time-Sharing System. Comm. ACM 17, (July 1974), p365.
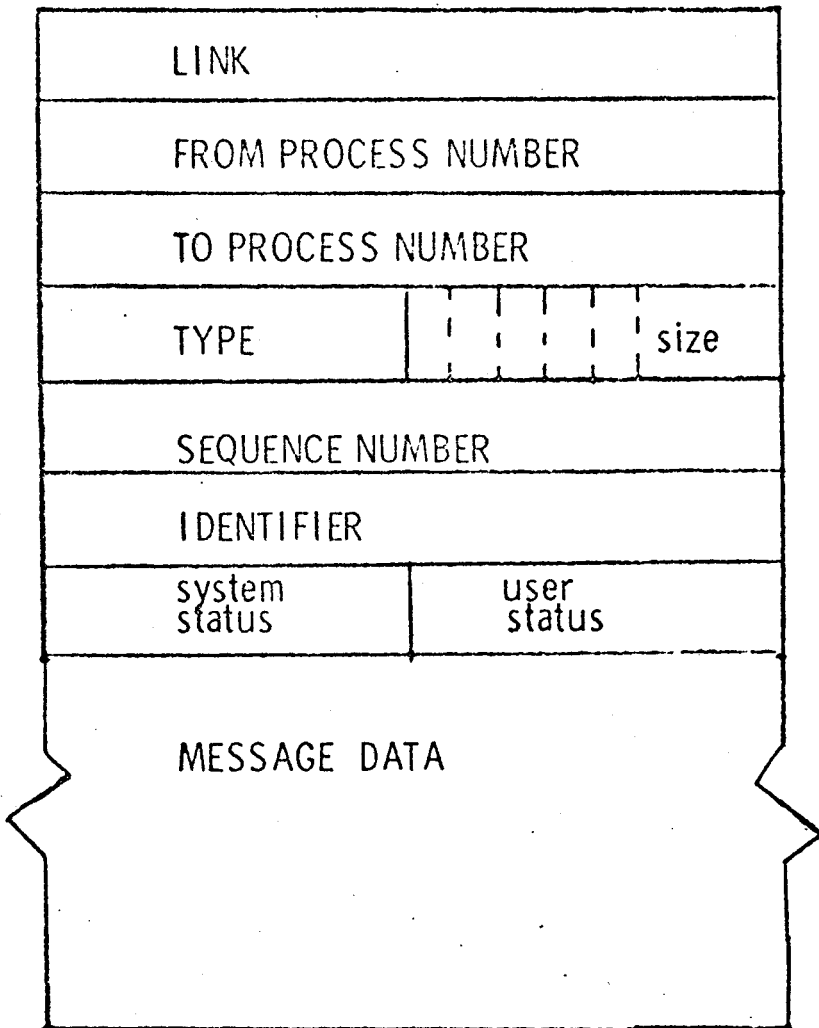
Figure 1 System Structure

FIGURE 2 Message format