# On Reliable and Extendible Operating Systems

## INTRODUCTION

A considerable amount of bitter experience in the design of
operating systems has been accumulated in the last few years, both
by the designers of the systems which are currently in use and by
those who have been forced to use them. As a result, many people
have been led to the conclusion that some radical changes must be
made, both in the way we think about the functions of operating
systems and in the way they are implemented. Of course, these are
not unrelated topics, but it is often convenient to organize ideas
around the two axes of function and implementation.

This paper is concerned with an effort to create more flexible and
more reliable operating systems built around a very powerful and
general protection mechanism. The mechanism is introduced at a
low level and is then used to construct the rest of the system,
which thus derives the same advantages from its existence as do
user programs operating with the system. The entire design is
based on two central ideas. The first of these is that an opera-
ting system should be constructed in layers, each one of which
creates a different and hopefully more convenient environment in
which the next higher layer can function. In the lower layers a
bare machine provided by the hardware manufacturer is converted
into a large number of *user machines* which are given access to
common resources such as processor time and storage space in a
controlled manner. In the higher layers these user machines are
made easy for programs and users at terminals to operate. Thus,
as we rise through the layers, we observe two trends.
1   The consequences of an error become less severe
2   The facilities provided become more elaborate
At the lower levels we wish to press the analogy with the hardware
machine very strongly; *where the integrity of the entire system is
concerned, the operations provided should be as primitive as
possible.* This is not to say that the operations should not be

complete, but that they need not be convenient. They are to be
regarded in the same light as the instructions of a central
processor. Each operation may in itself do very little; we require
only that the entire collection should be powerful enough to permit
more convenient operations to be programmed.

The main reason for this dogma is clear enough; simple operations
are more likely to work than complex ones and, if failures are to
occur, it is very much preferable that they should hurt only one
user rather than the entire community. We therefore admit increas-
ing complexity in higher layers, until the user at his terminal may
find himself invoking extremely elaborate procedures. The price to
be paid for low level simplicity is also clear: additional time to
interpret many simple operations and storage to maintain multiple
representations of essentially the same information. We shall
return to these points below. It is important to note that users
of the system other than the designers need not suffer any added
inconvenience from its adherence to the dogma, since the designers
can very well supply, at a higher level, programs that simulate
the action of the powerful low level operations to which users may
be accustomed. Users do, in fact, profit from the fact that a
different set of operations can be programmed if the ones provided
by the designer prove unsatisfactory. This point also will receive
further attention.

The increased reliability that we hope to obtain from an applica-
tion of the above ideas has two sources. In the first place,
careful specification of an orderly hierarchy of operations will
clarify what is going on in the system and make it easier to
understand. This is the familiar idea of modularity. Equally
important, however, is a second and less familiar point, the other
pillar of our system design, which might loosely be called enforced
modularity. It is this: if interactions between layers or modules
can be forced to take place through defined paths only, then the
integrity of one layer can be assured regardless of the deviations
of a higher one. The requirement that no possible action of a
higher layer, whether accidental or malicious, can affect the
functioning of a lower one is a strong one. In general, hardware
assistance will be required to achieve this goal, although in some
cases the discipline imposed by a language such as ALGOL, together
with suitable checks on the validity of subscripts, may suffice.
The reward is that errors can be localized very precisely. No
longer does the introduction of a new piece of code cast doubt on
the functioning of the entire system, since it can affect only its

own and higher layers.


CONTEXT FOR EXAMPLES


In order to keep the remainder of the paper grounded in reality,
most of the ideas will be discussed with reference to a specific
system, within the context of which many of them were indeed
evolved.  This is the time-sharing system being developed for the
dual processor Control Data 6400 at the University of California,
Berkeley.  I present here a brief sketch of the relevant features
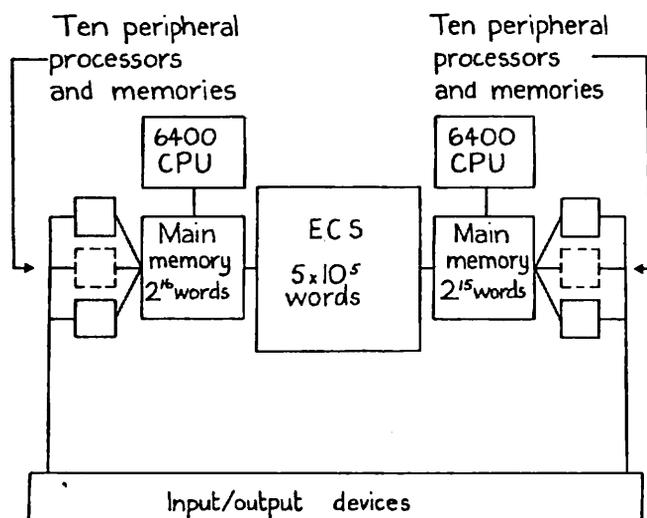of this system.



*Figure 1: Dual 6400 hardware configuration*


Four aspects of the hardware are important to us (see Figure 1).
1    The two processors have independent main core memories and
     communicate through the extended core storage (ECS), which has
     a latency of 4 microseconds and is capable of transferring 10
     words of 60 bits each in one microsecond.  This means that a
     program of 10 000 words, or about 30 000 instructions, can be
     swapped into core in one millisecond.  The ECS is regarded as
     the normal repository for active data.  Only one user program
     at a time is held in each main core memory.
2    Each processor has very simple memory protection and mapping
     consisting of a relocation register and a bounds register.
3    The entire state of a processor, all the information not in
     memory that is required to define a running program, can be
     switched in 2 microseconds.

4    Input/output is handled by ten peripheral processors for each
     central processor.  They all run independently, each with its
     own 4K x 12-bit memory.  All can access main memory, but not
     ECS, all can switch the state of the central processor and all
     have exactly the same access to the input/output devices of
     the system.  There are no interrupts.

The software is organized into a basic system, called the ECS
system, which runs on the bare machine and consists of a single
(lowest) layer for protection purposes and any number of modules
which may form higher layers.  The ECS system implements a small
number of basic *types* of *objects*, each one of which can be named
and referred to independently (see Figure 2).

```
File
Process
Event channel
Capability list (C-list)
Operation
Class code
Allocation block
```

*Figure 2: Types of objects in the ECS system*

Data in the system is stored in files.  Each file is an ordered
sequence of words that are numbered starting at zero.  Operations
exist to address a block of words in a file and transfer data
between the block and memory.

A process is a vehicle for the execution of programs, a logical
processor which shares the physical processor with other programs.
It consists of a machine state, some resources which it expends in
doing work and some additional structure to describe its memory
and its access to objects.  The details of this structure are the
subject of later sections of this paper.

Processes communicate through shared files or through event
channels, which are first-in, first-out queues of 1-word messages
or *events*.  A process may try to read an event from one of a list
of event channels, in which case it will be blocked from further
execution until an event is sent to one of the channels by another
process.  Many processes may be blocked on the same event channel
and many events may be sent to one channel before any process
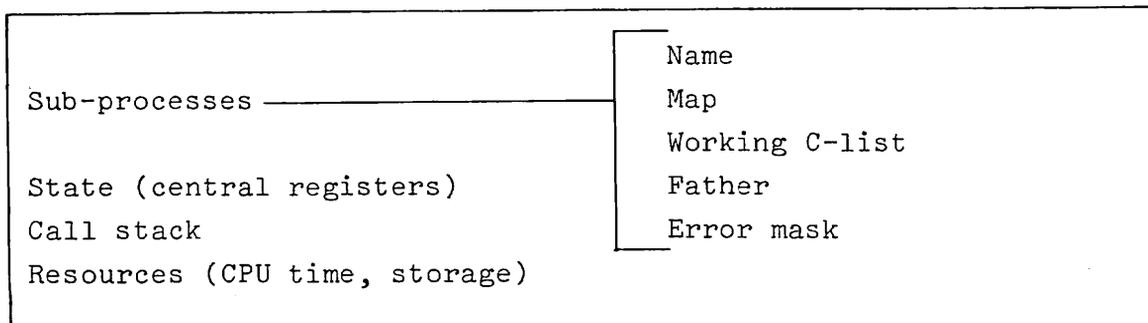comes to read them.

```
┌──────────────────────────────────────────────────────────┐
│                                   ┌─ Name                 │
│  Sub-processes ───────────────────┤  Map                  │
│                                   │  Working C-list       │
│  State (central registers)        │  Father               │
│  Call stack                       └─ Error mask           │
│  Resources (CPU time, storage)                            │
└──────────────────────────────────────────────────────────┘
```

*Figure 3: Components of a process*

Allocation blocks are used to control and account for the expenditure of system resources which, in the ECS system, are only storage space in ECS and CPU time. Every object in the system is owned by an allocation block, which provides the resources for the space the object takes up and accumulates charges for the word-seconds of storage expended in keeping the object in existence. Allocation blocks also allow all the objects in the system to be found, since they form a tree structure rooted in a single block belonging to the system.

The remaining types of objects in the ECS system are closely related to the subject matter of this paper; we now turn to consider them in more detail.

NAMES AND ACCESS RIGHTS

All objects in the system are named by *capabilities*, which are kept in capability lists or C-lists. These lists are like the memory of a rather peculiar 2-address machine, in the sense that operations exist to zeroize C-list entries and to copy capabilities from entry $i$ of C-list $A$ to entry $j$ of C-list $B$. In addition, any operation that creates an object deposits a capability for it in a designated C-list entry. The function of a capability is two-fold.
1    It names an object in the system.
2    It establishes a right to do certain things with the object.
At any given time a process has a single working C-list, W. A capability is referenced by specifying an entry in W; the $i$th entry will be referred to as W[$i$]. Since C-lists are objects that can themselves be named by capabilities, it is possible to specify capabilities in more complex ways; for example, W[$i$][$j$] would be the $j$th entry of the C-list named by the capability in the $i$th

entry of W.   In the interests of simplicity, however, all
capabilities passed to operations in the ECS system must be in W
and they are specified by integers that refer to entries of W.

In this rather obvious way a capability, which is the protected
name of an object (that is, it cannot be altered by a program),
itself acquires an unprotected name.   The integrity of the
protection system is preserved because only capabilities in W can
be so named.   The fact that a capability is in W is thus construed
as prima facie evidence that the process has the right to access
the object that it names.   From this foundation a complex directed
graph of C-lists may be built up which provides a great deal of
flexibility and convenience in the manipulation of access rights,
although it is still limited in certain important respects which
we shall explore later.

A capability is actually implemented as two 60-bit words (see
Figure 4).   The type field is an integer between 1 and 7 if the
capability is for an object defined by the ECS system.   Other
values of the type are for user-created capabilities, which are
discussed below.   The MOT index points to an entry in the Master
Object Table, which in turn indicates where to find the object,
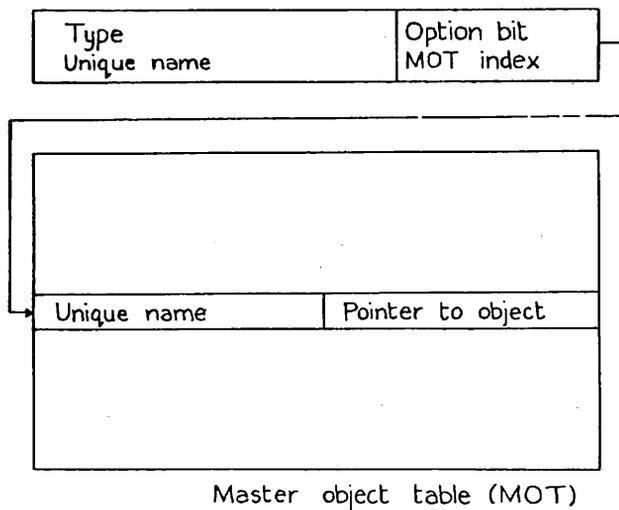that is, gives its address in ECS.

| Type | Option bit |
| Unique name | MOT index |

|  |  |
| Unique name | Pointer to object |
|  |  |

Master object table (MOT)

*Figure 4: The structure of a capability*

The unique name, guaranteed to be different for each object in the
system, must be the same in the capability and in the MOT entry.
This requirement makes it possible to destroy an object and re-use
its MOT entry without finding and destroying all the capabilities

for it.  If we could not do this, it would be necessary to restrict severely the copying of capabilities and to keep an expensive and error-prone list of back-pointers from each object to its capabilities.  Two additional benefits are obtained from the unique name organization of the MOT.

1   If an object is moved in ECS, only the pointer to it in the MOT needs to be updated, since all references to the object are made by indirection through the MOT.  Since the system's storage allocation strategy is a first-fit search of free blocks, followed by a compacting of free space if no block is big enough, it is essential to be able to move objects.

2   If some damage is accidentally done to a capability, by hardware failure or a bug in the ECS system, it is extremely unlikely that the damaged capability can be used improperly, since the chance that the unique name will match with the one in the MOT is small.  It is worthwhile to note that a slightly modified version of this scheme, in which the MOT index is dispensed with and the object is found by association on the unique name, is also possible, although significantly more expensive to use.

The option bits field of a capability serves as an extension of the type.  Each bit of the field should be thought of as authorizing, if it is set, some kind of operation on the object.  A file capability, for example, has option bits authorizing reading, writing, deletion and various more obscure functions.  The operation that copies a capability allows any of the option bits to be turned off, so that weaker capabilities can easily be made from a given one in a systematic way and without a host of special operations.  The interpretation of option bits is also systematized by the system's treatment of operations, to which we now turn.

OPERATIONS

Operations can be thought of as the instruction set of a user machine created by the system, or as the means by which a program examines and modifies its environment.  Viewing them in these ways, we want them to have the following features.

1   Operations must be able to be handed out selectively, so that the powers exercised by a program can be controlled.

2   The mapping between the names that a program uses to specify operations and the operations themselves must be able to be

The header at top left is the running header.

changed, so that it is possible to run a program in an
"envelope" and alter the meaning of its references to the
outside world.

3  Users must be able to create new operations that behave in
exactly the same way as the operations originally provided by
the ECS system must.

4  A systematic scheme must exist to handle error conditions that
may arise during the attempted execution of an operation.

The first two points are dealt with by treating operations as
objects for which capabilities are required.  This means that a
process can call on only those operations that it finds in its
working C-list, W.  Furthermore, since operations are identified
only by indices in W, the meaning of a call on operation 3, say,
can easily be changed by changing the contents of W[3].  When a
program starts to run, it expects to find at pre-arranged locations
in W the operations that it needs in order to function.  All of its
communication with the outside world is therefore determined by the
initial contents of W.

We now proceed to consider the internal structure of an operation
in more detail.  An operation is a sequence of *orders* which are
numbered starting at 1.  Each order consists of an *action* and a
parameter specification list, which is a sequence of *parameter
specifications*.  The action tells what to do; it is either an ECS
system action or a user-defined action (discussed below).  The
parameter specification list describes the parameters that are
expected.  Each parameter may be:

1  A data word, which is simply a 60-bit number.

2  A capability, in which case the parameter specification
specifies the type and the option bits that must be set in
the actual parameter.

When the operation is called (see Figure 5), a list of prototype
actual parameters must be supplied.  Each one is a number.  If the
corresponding parameter specification for order 1 calls for a data
word, the number itself becomes the actual parameter; if the
parameter specification calls for a capability, the number is
interpreted as an index in W and the capability W[$i$] becomes the
actual parameter, provided that the type is correct and all the
option bits demanded by the parameter specification are set in
W[$i$].

Given an operation, it is possible to create from it a new
operation in which some of the parameter specifications are

|   | Type | Option | *Unique* name |
|---|------|--------|------|
| *1* | File | 1010 | 6341 |
| *2* | File | 0011 | 2533 |
| *3* | Op | 0000 | 6677 |
| *4* | Process | 0001 | 2431 |
|   |   |   |   |
|   | . |   |   |
|   | . |   |   |
|   | . |   |   |

*Working C-list*

Call 3,    Operation 6677    *Actual parameters*

*Prototype actual parameters*

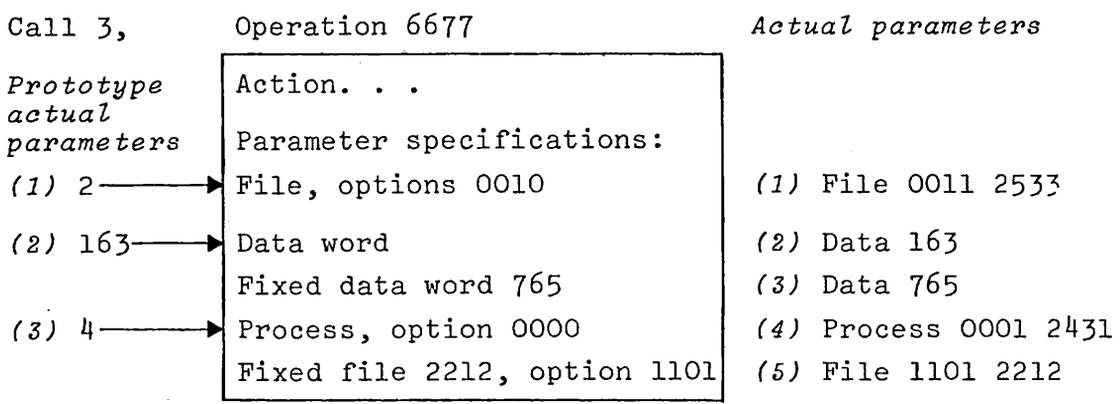| | Action. . . | |
|---|---|---|
| *(1)* 2 ——→ | Parameter specifications:<br>File, options 0010 | *(1)* File 0011 2533 |
| *(2)* 163 ——→ | Data word | *(2)* Data 163 |
| | Fixed data word 765 | *(3)* Data 765 |
| *(3)* 4 ——→ | Process, option 0000 | *(4)* Process 0001 2431 |
| | Fixed file 2212, option 1101 | *(5)* File 1101 2212 |

*Figure 5: Calling an operation: constructing the actual parameter list*

converted into *fixed parameters*; that is, the actual parameters for the action are built into the operation and are no longer supplied in the prototype actual parameter list. In this way a general operation may be specialized in various directions without the addition of any significant overhead to a call.

An action may *return* a single numeric value after it has completed successfully. Of course, the caller can pass it files and C-lists in which it can return either data or capabilities of arbitrary complexity. It may also *fail* if it meets with some circumstances beyond its competence. In this case it has two options: it may return with an *error*, which is handled by mechanisms described later; alternatively, it may take a *failure return*. The subsequent action of the system depends on the order structure of the operation. When the call on the operation is made, the parameter specification list of order 1 is used to interpret the arguments and the action of order 1 is executed.

The *order level*, $i$, of the call is set to 1. If the action takes

a failure return, the system re-examines the operation to see if
it has order $i + 1$.  If it hasn't, a failure return is made to the
caller.  If it has, $i$ is increased by 1 and order $i$ of the operat-
ion is called.

The rationale behind this rather elaborate mechanism is to allow
relatively simple operations to be supported by more complex ones.
Suppose that $A$ is a simple and cheap operation on a file and it
fails under certain, hopefully rare, circumstances.  For example,
$A$ might read from the file and might fail if no data is present.
Now operation $B$ may be devised to recover from the failure of $A$;
it might attempt to obtain the missing data from disk storage.
From $A$ and $B$ we make the 2-order operation $C = (A,B)$.  A call of
$C$ now costs no more than a call of $A$ if the data is present in the
file.  If it is not, $B$ must be called at considerable extra cost,
but this is hopefully an infrequent occurrence.  The alternative
approach is to call $B$ and have it in turn call $A$.  This is quite
unsatisfactory if we are thinking in terms of a system that may
eventually have many layers, since it requires passage through
every layer to reach the most basic, bottom layer.  Such a design
is acceptable if each layer expands the power of the operation, so
that a great deal more work is normally done by a call on layer 2
than by a call on layer 1; it is not acceptable, however, when the
higher layers are present to deal with unusual conditions, without
normally adding anything to the work accomplished.


## USER-DEFINED OPERATIONS AND PROTECTION

The last section has suggested two ways to look at an operation.
1   As a machine instruction in an extended or user machine
2   As a means of communicating with the world outside a program
A third analogy which is inevitably suggested is with an ordinary
subroutine call.  The crucial difference between a call on an
operation and a subroutine call is that the former involves a
change in the capabilities accessible to the process, that is, in
the working C-list.  This is obviously the case when the operation
is one defined by the ECS system, since the code that implements
the operation is then running entirely outside of the system's
protection structure.  For a user-defined operation, a more formal
mechanism must exist within the overall protection structure for
specifying how the capabilities of the process change when the
operation is called.

To this end, some additional structure is defined for a process (see Figure 3). In particular, a new entity called a *sub-process* is introduced. Associated with each sub-process is a working C-list and a map which defines the memory of the sub-process. At any given instant the process is executing in one *active* sub-process, or in the ECS system itself, and consequently has the working C-list and memory of that sub-process. When the active sub-process changes, the memory that the process addresses and the working C-list also change, so that the process finds itself in a different environment.

Because a change in the active sub-process, that is, a transfer of control from one sub-process to another, implies a change in capabilities, there must be some means of controlling the ways in which such transfers are allowed to take place. This is done as follows. A sub-process can be called only by calling on an operation that has that sub-process as its action; the means for constructing such operations are discussed below. The call proceeds as follows:
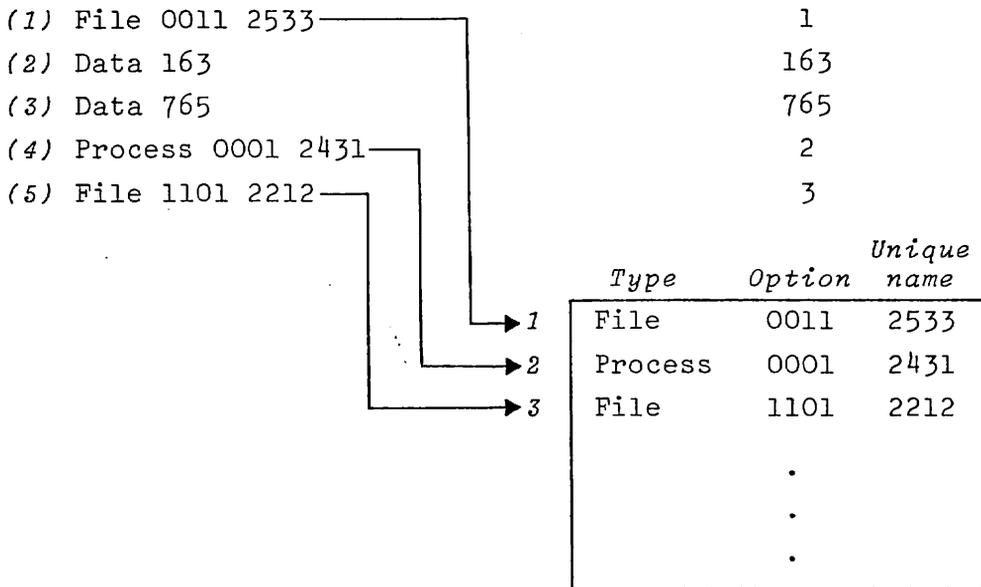
1    Compute the actual parameter list. The rules for doing this are described above; Figure 5 gives an example.

2    Copy a representation of the actual parameter list into a fixed place in the memory of the sub-process being called. Figure 6 illustrates this procedure. An actual parameter that is a data item, (2) in Figure 6, for example, is represented by its index in W. An actual parameter that is a capability, (4) in Figure 6, for example, is copied into the working C-list of the sub-process being called and is represented by its index in that C-list. Note that the representation of the actual parameter list could be used as the prototype actual parameter list for another call on the operation.

3    Start executing the sub-process being called at a fixed location called the *entry point*.

This calling mechanism has been designed with some care to have the following features:

1    It is possible to control who can call a sub-process by controlling creation and distribution of operations with the necessary action.

2    Calls can enter the sub-process only at a single point.

3    The called sub-process need not have either fewer or more capabilities than the caller. Any capabilities needed by the caller can be passed as parameters. If capabilities are to be returned, the caller can pass a C-list into which the called

| Actual parameters<br>(from Figure 5) | | Representation of actual<br>parameters in the called<br>sub-process |
|---|---|---|
| *(1)* File 0011 2533 | | 1 |
| *(2)* Data 163 | | 163 |
| *(3)* Data 765 | | 765 |
| *(4)* Process 0001 2431 | | 2 |
| *(5)* File 1101 2212 | | 3 |

|   | Type | Option | Unique<br>name |
|---|---|---|---|
| 1 | File | 0011 | 2533 |
| 2 | Process | 0001 | 2431 |
| 3 | File | 1101 | 2212 |
|   | . | | |
|   | . | | |
|   | . | | |

*Working C-list for called*
*sub-process*

*Figure 6: Calling an operation: copying the actual parameter list*
sub-process deposits the new capabilities.

4    An operation that calls a sub-process is used exactly like one
that calls the ECS system or any other sub-process.  It is
therefore very easy to run a program and intercept some or all
of its calls on the ECS system or on the sub-processes.

A call operation implies a return.  Since we do not want any
unnecessary restrictions on the depth to which calls can take
place or on recursion calls, a stack is used to keep track of
return points.  It is referred to as the *call stack* and maintained
by the ECS system.

ACCESS TO SUB-PROCESSES

As the careful reader may have noticed, we have refrained from
saying that a sub-process is an object in the system.  There are
two reasons for this.  Firstly, a sub-process does not exist
independently of the process that contains it.  Secondly, and more
importantly, we wish to have a means of referring to similar sub-
processes in different processes by the same name, so that the

same user-defined operations can be shared by a number of processes.
To achieve this goal, we introduce the last type of ECS system
object.  It is called a *class code* and consists simply of a 60-bit
number.  This number is divided into two parts, called the
*permanent* and *temporary* parts.  There are two operations on class
codes: ·

1   Create a class code with a new permanent part, never seen
    before, and a temporary part equal to zero.

2   Set the temporary part of a class code.  This operation
    requires one of the option bits to be set.  It is therefore
    possible to fix the temporary part also, simply by turning
    off this bit.

A class code can be thought of as a rather flexible means for
authorizing various things to be done without the need for a large
number of individual capabilities.  A sub-process, for example, is
named by a class code.  This means that anyone with a capability
for the class code can create an operation that calls on the sub-
process with that name.  When the operation is created, a capabi-
lity for it is returned and can then be passed to other processes
like any other capability.  If the class code is passed along with
it, sub-processes with the same name can be created in several
processes and called with the same operation.  Since the permanent
part of a class code is unique, there is no chance that independ-
ently created operations will name the same sub-process. Further-
more, it is easy to create a sub-process that cannot be called by
any existing operations, simply by creating a new class code and
then using it to name the sub-process.

The most useful application of the scheme is in connection with
sub-systems consisting of a number of operations, various files
containing code and data and, perhaps, several sub-processes.  All
the necessary information can be gathered together into one C-list
and used to create copies of the sub-system in any number of
processes.

One other application for class codes is built into the ECS system:
they are used to authorize the creation of user-defined types of
capability.  Thus there is an operation that takes a class code
and a data word and creates a capability (see Figure 4) with the
type given by part of the class code and the second word given by
the data word.  Such capabilities will always produce errors if
given to ECS system operations, but they may be passed successfully
to user-defined operations.  In this way users can create their own

kinds of objects and take advantage of the ECS system facilities
for controlling access to them.


## ERRORS

A running program can cause errors in a variety of ways: by
violating the memory protection, executing undefined operation
codes, or making improper calls on operations.  Some orderly method
is required for translating these errors into calls on other
programs that will take responsibility for dealing with them.  In
order to do this, it is necessary to attach to each sub-process,
S, another one to which errors should be passed if S is unwilling
to handle them.  We call this other sub-process the *father* of S
and insist that the father relation define a tree structure on the
sub-processes of a process so that an error condition can be passed
from one sub-process to another along a path that eventually
terminates.

When an error occurs, it is identified by two integers called the
*error class* and the *error number*.  Every sub-process has a bit
string called its *error selection mask* (ESM), and is said to *accept*
an error if the bit of its ESM selected by the class is on.  When
an error occurs, a search is made for a sub-process, *A*, that
accepts it, starting with the sub-process in which it occurs and
proceeding along the path defined by the father relation.  The
root of the tree structure is assumed to accept all errors.  When
*A* is found, it is called with the error class and error number as
arguments.  The entire current state of things is thus preserved
for *A* to examine.  If it wants to patch things up and continue
execution, it can just return.  If it decides to abort the compu-
tation, it can force a return to some sub-process farther up the
call stack.


## MEMORY AND MAPS

It is fairly obvious that the right to access the memory addressed
by a program is similar in nature to the right to access a file.
Logically, it should therefore be controlled by a capability which
should be mentioned every time an access is made.  On a segmented
machine this is a very natural and satisfactory point of view.

Lacking segmentation hardware, however, we must adopt a variety of compromises. Further complications are introduced by the fact that many machines, including the 6400, do not have address spaces large enough to allow all sub-processes to talk about each other's memory, and certainly no satisfactory controls on access to it. This section is concerned with the rather ad hoc schemes adopted in the system under discussion to deal with this problem. It is included in the paper for two reasons:

1   To point out the salient problems to be faced by any system.
2   To show how unpleasant are the expedients to which unsuitable hardware may force us.

No attempt is made to discuss the problem in the abstract or with any generality and no argument is offered in favor of the devices described, except expediency on the available hardware.

The first problem is to find a representation of a process' memory when the process is not running and therefore is not in main memory. In order to avoid introducing a new kind of object and to facilitate the sharing of read-only code, a scheme suggested by hardware mapping mechanisms has been adopted. The memory of a sub-process is defined by a *map*, which consists of a list of entries of the form:

memory address *m*, file *f*, file address *a*, length *l*, read-only flag *r*.

The strict meaning of such an entry is this:

1   When the process is chosen to run next, *l* words starting with word *a* in file *f* are transferred into main memory starting at word *m* (relative to the location addressed by the hardware relocation register).
2   When the process stops running, the transfer is reversed if *r*=0. Otherwise, nothing is done.

We usually think of this in the terms suggested by Figure 7; a section of the file is made to correspond to a section of the addressable storage of the sub-process. The analogy breaks down if any modification is made to the file while the main memory copy produced by the map entry exists; this is extremely unfortunate if writeable data bases are to be shared, but it is unavoidable. The scheme works very well for sharing read-only programs and data, however. Multiple copies need exist only in central memory, which is not a precious resource since only one process is allowed to reside there at a time because of the very high swap rate.

The second problem is how one sub-process can access the memory of another one. A straightforward solution is to confine data sharing
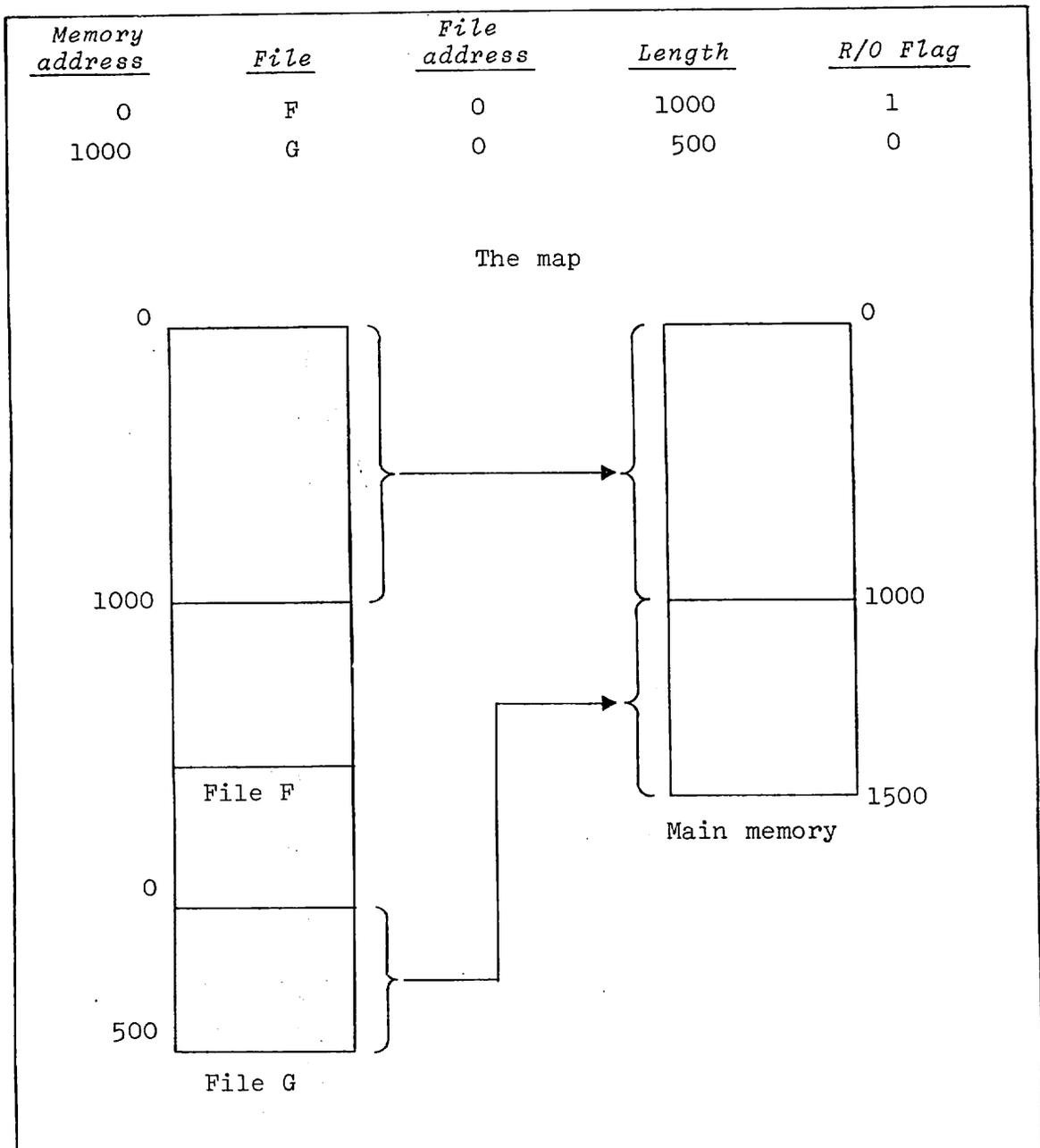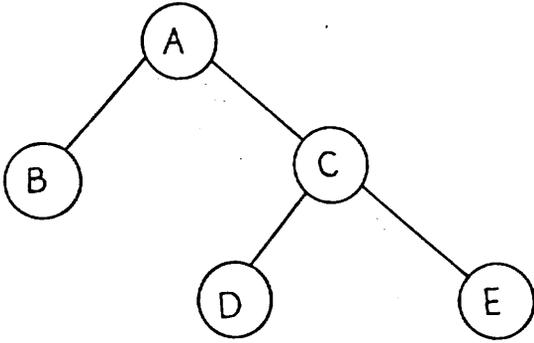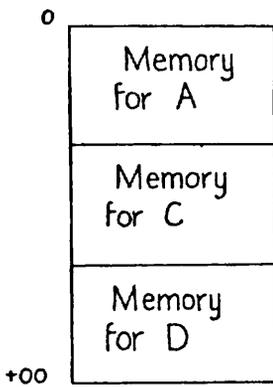
| Memory address | File | File address | Length | R/O Flag |
|---|---|---|---|---|
| 0 | F | 0 | 1000 | 1 |
| 1000 | G | 0 | 500 | 0 |

The map



*Figure 7: Operation of the map*

to files, but cursory examination reveals that this is extremely
inconvenient, since the relationship between memory and files
established by the map is quite indirect.  On the other hand,
appending the address space of the caller to the called program is
also unacceptable, both because such a sledgehammer approach
negates the selectivity of the rest of the protection system and
because the sum of the lengths of the two address spaces may
easily exceed the available central memory.  A restricted solution
is to allow a sub-process to append to its address space the spaces
of all its descendants along a single path to a leaf of the tree
(see Figure 8).  The map for this extended address space is called
a *full map* and is obtained by concatenating the maps of the sub-
processes involved.  It is constructed automatically whenever a

Sub-process tree



One possible full map for A

*Figure 8: Full maps*

sub-process calls one of its ancestors, since the path to be used
is then uniquely defined. With this scheme we have achieved:

1   Convenient addressing across sub-process boundaries for all
    the sub-processes on a path held together by father pointers.

2   Restriction of the total address space size for all the
    sub-processes on any such path to the total available main
    memory.

Note that, when constructing a full map, it is never necessary to
move anything, since when *D* is swapped in, say, the space required
by *A* and *C* is known whether or not they are swapped in also.

## SYSTEM EXTENDABILITY

In two concluding sections we present some thoughts on the general
properties that we should expect an operating system to possess, if
it is to be a firm foundation for the construction and operation of
software, including itself, and on methods for realizing these
properties.  The major purpose of the detailed description in the
preceding sections was to provide a concrete illustration of these
ideas.

If a system is to evolve to meet changing requirements and if it is
to be flexible enough to permit modularization without serious
losses of efficiency, it must have a basic structure that allows
extensions not only from a basic system but also from some complex
configuration that has been reached by several prior stages of
evolution.  In other words, the extension process must not exhaust
the facilities required for further extensions.  The system must be
completely open-ended, so that additional machinery can be attached
at any point.

Secondly, the relations between a module and its environment must
be freely re-definable.  This means that any module, provided its
references to the outside world are documented, can be run inside
an "envelope" that intercepts all of these references and re-
interprets them at will.  In order to ensure that external refer-
ences are documented, it must be normal, and indeed compulsory,
practice to specify each module's environment precisely and
independently of the module's internal structure.  This requirement
is satisfied by the use of C-lists to define the outside world to
a sub-process.

Thirdly, it must be possible to introduce several layers of re-
interpretation around a module economically, without forcing all
of its external references to suffer re-interpretation by each
layer.  In other words, a capability for extension by exception is
required.  Furthermore, the system's basic mechanisms for naming
and for calls must leave room for a number of higher level sub-
systems to make their mark, rather than forcing each new sub-system
to create and maintain its own inventory of objects.

To summarize, a usefully extendable system must be open-ended, must
allow a sub-system to be isolated in an envelope and must encourage
economical re-use of existing constructs.  Such a system has some
chance of providing a satisfactory toolkit for others to use.

SYSTEM RELIABILITY

It is even more important that a useful system should be a
functioning one.  Since we do not know how to guarantee the
correctness of a large collection of interacting components, we
must be able to break systems up into units in such a way that:
1   Each unit is simple enough to be fully debugged.
2   Each unit interacts with only a few other units.
If this division is to inspire confidence, it must be enforced.
It is not feasible to depend on every contributor for good will
and a full understanding of the rules for intercourse with others.
Hence a complete and precise protection system is needed.

A great deal of flexibility is required in the manipulation of
access rights; otherwise the protection facilities will prove so
cumbersome to use that they will quickly be abandoned in favor of
large, monolithic designs.  It must become a pleasure to write
programs that safeguard themselves against the inroads of others,
or at the very least it must be automatic, almost as unavoidable
as the use of a higher level language.

Thirdly, the implementation must be fail-fast: it should detect a
potential malfunction as early as possible so that corrective
action can be taken.  The use of unique names with pointers,
redundant data structures, parity bits and checksums are all
valuable devices for warning of impending disaster.  On the other
hand, elaborate and fragile pointer structures or allocation tables
that cannot be checked or reconstructed from the devices being
allocated, are likely to cost more than they are worth.

Most important, perhaps, is a general acceptance of the fact that
a flexible and reliable system will exact its price.  Under ideal
circumstances, the price will be paid in careful design and modest
amounts of special hardware to facilitate the basic operations of
the system.  More likely, though, is sizable amounts of software
overhead to make up for basic deficiencies in the machine.  Beyond
a certain point, admittedly not often reached, there is little that
can be done about this overhead.  It can be minimized by keeping
the goals of a system within reasonable bounds, but tends to be
increased by the final consideration in reliability.  This, of
course, is simplicity.

Figure 9 lists the operations of the 6400 ECS system.  There are
about 50 of them and few are implemented by more than a couple of

| | |
|---|---|
| Create file | Create sub-process |
| Delete file | Destroy sub-process |
| Create file block | Return |
| Delete file block | Failure return |
| Read shape | Jump return |
| Test for file block | Set map entry |
| Move back | Change map entry |
| Read from file | Display map entry |
| Write on file | Set ESM |
| | Set program counter |
| Create event channel | |
| Delete event channel | Create operation of order 1 |
| Send event | Fix PS to data |
| Read event | Fix PS to capability |
| | Copy operation |
| Create C-list | Add order to operation |
| Delete C-list | Delete operation |
| Display capability from C-list | |
| Display capability from W | Create allocation block |
| Copy capability and decrease options | Delete allocation block |
| | Transfer funds |
| Copy capability into W | Create capability for first object owned by block |
| Copy capability out of W | |
| Create process | Create new class code |
| Delete process | Change temporary part |
| Display state | |
| Send interrupt | Save registers |
| | Restore registers |
| | Change unique name of object |

*Figure 9: ECS system operations*

hundred instructions, most by fewer.  To convert the system they
define into one suitable for a general user will take several times
as much code at higher levels, but it rests on a secure foundation.