# EXHIBIT 13

Docket No.    375246US91RX

# IN THE UNITED STATES PATENT AND TRADEMARK OFFICE

IN RE REEXAMINATION OF: Richard Michael NEMES

| | | |
|---|---|---|
| CONTROL NO: | 90/011,426 | GAU:    3992 |
| FILED: | January 10, 2011 | EXAMINER:  ANDREW NALVEN |
| FOR: | METHODS AND APPARATUS FOR INFORMATION STORAGE AND RETRIEVAL USING A HASHING TECHNIQUE WITH EXTERNAL CHAINING AND ON-THE-FLY REMOVAL OF EXPIRED DATA | |

# INFORMATION DISCLOSURE STATEMENT UNDER 37 CFR 1.555

COMMISSIONER FOR PATENTS
ALEXANDRIA, VIRGINIA 22313

SIR:
Patent holder(s) wishes to disclose the following information.

## REFERENCES

■ Patent holder(s) wishes to make of record the reference(s) listed on the attached form PTO-1449 and/or accompanying documents from a corresponding foreign application. Copies of the listed reference(s) are attached, where required, as are either statements of relevancy or any readily available partial or full English translations of pertinent portions of any non-English language reference(s).

☐ Credit card payment is being made online (if electronically filed), or is attached hereto (if paper filed), in the amount required under 37 CFR §1.17(p).

## RELATED CASES

☐ Attached is a list of patent holder's pending application(s), published application(s) or issued patent(s) which may be related to the present application. In accordance with the waiver of 37 CFR 1.98 dated September 21, 2004, copies of the cited pending applications are not provided. Cited published and/or issued patents, if any, are listed on the attached PTO form 1449.

☐ Credit card payment is being made online (if electronically filed), or is attached hereto (if paper filed),  in the amount required under 37 CFR §1.17(p).

Respectfully submitted,

OBLON, SPIVAK, McCLELLAND,
MAIER & NEUSTADT, L.L.P.

Scott A. McKeown
Registration No. 42,866

Customer Number

# 22850

Tel. (703) 413-3000
Fax. (703) 413-2220
(OSMMN 02/10)

| Form PTO 1449 (Modified) | U.S. DEPARTMENT OF COMMERCE PATENT AND TRADEMARK OFFICE | ATTY DOCKET NO. 375246US91RX | | CONTROL NO. 90/011,426 |
|---|---|---|---|---|
| LIST OF REFERENCES CITED BY APPLICANT | | INVENTOR(S) Richard Michael NEMES | | |
| | | FILING DATE January 10, 2011 | | GROUP 3992 |

**U.S. PATENT DOCUMENTS**

| EXAMINER INITIAL | | DOCUMENT NUMBER | DATE | NAME | CLASS | SUB CLASS | FILING DATE IF APPROPRIATE |
|---|---|---|---|---|---|---|---|
| | AA | | | | | | |
| | AB | | | | | | |
| | AC | | | | | | |
| | AD | | | | | | |
| | AE | | | | | | |
| | AF | | | | | | |
| | AG | | | | | | |
| | AH | | | | | | |
| | AI | | | | | | |
| | AJ | | | | | | |
| | AK | | | | | | |
| | AL | | | | | | |
| | AM | | | | | | |
| | AN | | | | | | |

**FOREIGN PATENT DOCUMENTS**

| | | DOCUMENT NUMBER | DATE | COUNTRY | TRANSLATION | |
|---|---|---|---|---|---|---|
| | | | | | YES | NO |
| | AO | | | | | |
| | AP | | | | | |
| | AQ | | | | | |
| | AR | | | | | |
| | AS | | | | | |
| | AT | | | | | |
| | AU | | | | | |
| | AV | | | | | |

**OTHER REFERENCES (Including Author, Title, Date, Pertinent Pages, etc.)**

| | | |
|---|---|---|
| | AW | BAWDEN, Alan, et al., LISP Machine Progress Report - Massachusetts Institute of Technology Artificial Intelligence Laboratory - Memo No. 444; August 1977; 29 pgs. |
| | AX | Copyright 1995 by Bao Phan, et al.; Key Management Engine for BSD; 09/28/1995; Pgs. 1-29; (DEF00007942-DEF00007970) |
| | AY | Copyright 1995 by Bao Phan, et al.; Declarations and Definitions for Key Engine for BSD; 09/28/1995; Pgs. 1-4; (DEF00007971-DEF00007974) |
| | AZ | |

| Examiner | | Date Considered |
|---|---|---|

*Examiner: Initial if reference is considered, whether or not citation is in conformance with MPEP 609; Draw line through citation if not in conformance and not considered. Include copy of this form with next communication to applicant.

MASSACHUSETTS INSTITUTE OF TECHNOLOGY
ARTIFICIAL INTELLIGENCE LABORATORY

Memo No. 444

August 1977

LISP Machine Progress Report

by the Lisp Machine Group

Alan Bawden
Richard Greenblatt
Jack Holloway
Thomas Knight
David Moon
Daniel Weinreb

## ABSTRACT

This informal paper introduces the LISP Machine, describes the goals and current
status of the project, and explicates some of the key ideas.  It covers the LISP
machine implementation, LISP as a system language, input/output, representation
of data, representation of programs, control structures, storage organization,
garbage collection, the editor, and the current status of the work.

## INTRODUCTION:

The LISP Machine is a new computer system designed to provide a high performance and economical implementation of the LISP programming language.

The LISP language is used widely in the artificial intelligence research community, and is rapidly gaining adherents outside this group. Most serious LISP usage has historically been on the DEC PDP-10 computer, and both "major" implementations (InterLisp at BBN/XEROX and Maclisp at M.I.T.) were originally done on the PDP-10. Our personal experience has largely been with the Maclisp dialect of LISP, which was originally written in 1965.

Over the years, dramatic changes have taken place in the Maclisp implementation. At a certain point, however, modification and reimplementation of a language on a given machine can no longer efficiently gloss over basic problems in the architecture of the computer system. We, and many others, believe this is now the case on the PDP-10 and similar time-shared computer systems.

Time sharing was introduced when it became apparent that computers are easier to use in an interactive fashion than in a batch system, and that during an interactive session a user typically uses only a small fraction of the processor and memory available; often during much of the time his process is idle or waiting, and so the computer can be multiplexed among many users while giving each the impression that he is on his own machine.

However, in the Lisp community there has been a strong trend towards programs which are very highly interactive, very large, and use a good deal of computer time; such programs include advanced editors and debuggers, the MACSYMA system, and various programming assistants. When running programs such as these, which spend very significant amounts of time supporting user interactions, time sharing systems such as the PDP-10 run into increased difficulties. Not only is the processor incapable of providing either reasonable throughput or adequate response time for a reasonable number of users, but the competition for main memory results in large amounts of time being spent swapping pages in and out (a condition known as "thrashing"). Larger and larger processors and memory, and more and more complex operating systems, are required, with more than proportionally higher cost, and still the competition for memory remains a bottleneck. The programs are sufficiently large, and the interactions sufficiently frequent, that the usual time-sharing strategy of swapping the program out of memory while waiting for the user to interact, then swapping it back in when the user types something, cannot be successful because the swapping

cannot happen fast enough.

The Lisp Machine is a personal computer. Personal computing means that the processor and main memory are not time-division multiplexed, instead each person gets his own. The personal computation system consists of a pool of processors, each with its own main memory, and its own disk for swapping. When a user logs in, he is assigned a processor, and he has exclusive use of it for the duration of the session. When he logs out, the processor is returned to the pool, for the next person to use. This way, there is no competition from other users for memory; the pages the user is frequently referring to remain in core, and so swapping overhead is considerably reduced. Thus the Lisp Machine solves a basic problem of time sharing Lisp systems.

The user also gets a much higher degree of service from a Lisp machine than from a timesharing system, because he can use the full throughput capacity of the processor and the disk. Although these are quite inexpensive compared to those used in PDP-10 timesharing systems, they are comparable in speed. In fact, since disk access times are mainly limited by physical considerations, it often turns out that the disk used in a personal computer system is less expensive simply because of its smaller size, and has fully comparable throughput charactistics to the larger disk used by a timesharing system.

In a single-user machine, there is no penalty for interactiveness, since there are no competing users to steal a program's memory while it is waiting for its user to type. Thus the Lisp machine system, unlike time sharing systems, encourages highly interactive programs. It puts service to the user entirely ahead of efficiency considerations.

Another problem with the PDP-10 Lisp implementations is the small address space of the PDP-10 processor. Many Lisp systems, such as MACSYMA and Woods's LUNAR program, have difficulty running in an 18-bit address space. This problem is further compounded by the inefficiency of the information coding of compiled Lisp code; compilers for the PDP-10 produce only a limited subset of the large instruction set made available by the hardware, and usually make inefficient use of the addressing modes and fields provided. It is possible to design much more compact instruction sets for Lisp code. Future programs are likely to be quite a bit bigger; intelligent systems with natural language front ends may well be five or ten times the size of a PDP-10 address space.

The Lisp Machine has a 24 bit virtual address space and a compact instruction set, described later in this paper. Thus much larger programs may be used, without running into address space limitations. Since the instruction set is designed specifically for the Lisp language, the compiler is much simpler than

the PDP-10 compiler, providing faster and more reliable compilation.

The Lisp machine's compact size and simple hardware construction are likely to make it more reliable than other machines, such as the PDP-10; the prototype machine has had almost no hardware failures.

Much of the inspiration for the Lisp Machine project comes from the pioneering research into personal computing and display-oriented systems done by Xerox's Palo Alto Research Center.

## THE LISP MACHINE IMPLEMENTATION:

Each logged in user of the Lisp Machine system has a processor, a memory, a keyboard, a display, and a means of getting to the shared resources. Terminals, of course, are placed in offices and various rooms; ideally there would be one in every office. The processors, however, are all kept off in a machine room. Since they may need special environmental conditions, and often make noise and take up space, they are not welcome office companions. The number of processors is unrelated to the number of terminals, and may be smaller depending on economic circumstance.

The processor is implemented with a microprogrammed architecture. It is called the CONS Machine, designed by Tom Knight [CONS]. CONS is a very unspecialized machine with 32-bit data paths and 24-bit address paths. It has a large microcode memory (16K of 48-bit words) to accommodate the large amount of specialized microcode to support Lisp. It has hardware for extracting and depositing arbitrary fields in arbitrary registers, which substitutes for the specialized data paths found in conventional microprocessors. It does not have a cache, but does have a "pdl buffer" (a memory with hardware push-down pointer) which acts as a kind of cache for the stack, which is where most of the memory references go in Lisp.

Using a very unspecialized processor was found to be a good idea for several reasons. For one thing, it is faster, less expensive, and easier to debug. For another thing, it is much easier to microprogram, which allows us to write and debug the large amounts of microcode required to support a sophisticated Lisp system with high efficiency. It also makes feasible a compiler which generates microcode, allowing users to microcompile some of their

functions to increase performance.

The memory is typically 64k of core or semiconductor memory, and is
expandable to about 1 million words. The full virtual address space is stored on
a 16 million word disk and paged into core (or semiconductor) memory as required.
A given virtual address is always located at the same place on the disk. The
access time of the core memory is about 1 microsecond, and of the disk about 25
milliseconds. Additionally, there is an internal 1K buffer used for holding the
top of the stack (the PDL buffer) with a 200ns access time (see [CONS] for more
detail).

The display is a raster scan TV driven by a 1/4 Mbit memory, similar to
the TV display system now in use on the Artificial Intelligence Lab's PDP-10.
Characters are drawn entirely by software, and so any type or size of font can be
used, including variable width and several styles at the same time. One of the
advantages of having an unspecialized microinstruction processor such as CONS is
that one can implement a flexible terminal in software for less cost than an
inflexible, hardwired conventional terminal. The TV system is easily expanded to
support gray scale, high resolution, and color. This system has been shown to be
very useful for both character display and graphics.

The keyboard is the same type as is used on the Artificial Intelligence
Lab TV display system; it has several levels of control/shifting to facilitate
easy single-keystroke commands to programs such as the editor. The keyboard is
also equipped with a speaker for beeping, and a pointing device, usually a mouse
[MOUSE].

The shared resources are accessed through a 10 million bit/sec packet
switching network with completely distributed control. The shared resources are
to include a highly reliable file system implemented on a dedicated computer
equipped with state of the art disks and tapes, specialized I/O devices such as
high-quality hardcopy output, special-purpose processors, and connections to the
outside world (e.g. other computers in the building, and the ARPANET).

As in a time sharing system, the file system is shared between users.
Time sharing has pointed up many advantages of a shared file system, such as
common access to files, easy inter-user communication, centralized program
maintenance, centralized backup, etc. There are no personal disk packs to be
lost, dropped by users who are not competent as operators, or to be filled with
copies of old, superseded software.

The complete LISP Machine, including processor, memory, disk, terminal, and connection to the shared file system, is packaged in a single 19" logic cabinet, except for the disk which is free-standing. The complete machine would be likely to cost about $80,000 if commercially produced. Since this is a complete, fully-capable system (for one user at a time), it can substantially lower the cost of entry by new organizations into serious Artificial Intelligence work.

## LISP AS A SYSTEM LANGUAGE:

In the software of the Lisp Machine system, code is written in only two languages (or "levels"): Lisp, and CONS machine microcode. There is never any reason to hand-code macrocode, since it corresponds so closely with Lisp; anything one could write in macrocode could be more easily and clearly written in the corresponding Lisp. The READ, EVAL, and PRINT functions are completely written in Lisp, including their subfunctions (except that APPLY of compiled functions is in micro-code). This illustrates the ability to write "system" functions in Lisp.

In order to allow various low-level operations to be performed by Lisp code, a set of "sub-primitive" functions exist. Their names by convention begin with a "%", so as to point out that they are capable of performing unLispy operations which may result in meaningless pointers. These functions provide "machine level" capabilities, such as performing byte-deposits into memory. The compiler converts calls to these sub-primitives into single instructions rather than subroutine calls. Thus Lisp-coded low-level operations are just as efficient as they would be in machine language on a conventional machine.

In addition to sub-primitives, the ability to do system programming in Lisp depends on the Lisp machine's augmented array feature. There are several types of arrays, one of which is used to implement character strings. This makes it easy and efficient to manipulate strings either as a whole or character by character. An array can have a "leader", which is a little vector of extra information tacked on. The leader always contains Lisp objects while the array often contains characters or small packed numbers. The leader facilititates the use of arrays to represent various kinds of abstract object types. The presence in the language of both arrays and lists gives the programmer more control over data representation.

A traditional weakness of Lisp has been that functions have to take a fixed number of arguments. Various implementations have added kludges to allow variable numbers of arguments; these, however, tend either to slow down the function-calling mechanism, even when the feature is not used, or to force peculiar programming styles. Lisp-machine Lisp allows functions to have optional parameters with automatic user-controlled defaulting to an arbitrary expression in the case where a corresponding argument is not supplied. It is also possible to have a "rest" parameter, which is bound to a list of the arguments not bound to previous parameters. This is frequently important to simplify system programs and their interfaces.

A similar problem with Lisp function calling occurs when one wants to return more than one value. Traditionally one either returns a list or stores some of the values into global variables. In Lisp machine Lisp, there is a multiple-value-return feature which allows multiple values to be returned without going through either of the above subterfuges.

Lisp's functional orientation and encouragement of a programming style of small modules and uniform data structuring is appropriate for good system programming. The Lisp machine's micro-coded subroutine calling mechanism allows it to also be efficient.

Paging is handled entirely by the microcode, and is considered to be at a very low level (lower level than any kind of scheduling). Making the guts of the virtual memory invisible to all Lisp code and most microcode helps keep things simple. It would not be practical in a time sharing system, but in a one-user machine it is reasonable to put paging at the lowest level and forget about it, accepting the fact that sometimes the machine will be tied up waiting for the disk and unable to run any Lisp code.

Micro-coded functions can be called by Lisp code by the usual Lisp calling mechanism, and provision is made for micro-coded functions to call macro-coded functions. Thus there is a uniform calling convention throughout the entire system. This has the effect that uniform subroutine packages can be written, (for example the TV package, or the EDITOR package) which can be called by any other program. (A similar capability is provided by Multics, but not by ITS nor TENEX).

Many of the capabilities which system programmers write over and over again in an ad hoc way are built into the Lisp language, and are sufficiently good in their Lisp-provided form that it usually is not necessary to waste time worrying about how to implement better ones. These include symbol tables, storage management, both fixed and flexible data structures, function-calling,

and an interactive user interface.

Our experience has been that we can design, code, and debug new features much faster in Lisp-machine programs than in PDP-10 programs, whether they are written in assembler language or in traditional "higher-level" languages.


## INPUT/OUTPUT:

### Low level:

The Lisp Machine processor (CONS) has two busses used for accessing external devices: the "XBUS", and the "UNIBUS". The XBUS is 32 bits wide, and is used for the disk and for main memory. The UNIBUS is a standard PDP-11 16-bit bus, used for various I/O devices. It allows commonly available PDP-11 compatible devices to be easily attached to the Lisp Machine.

Input/output software is essentially all written in Lisp; the only functions provided by the microcode are %UNIBUS-READ and %UNIBUS-WRITE, which know the offset of the UNIBUS in physical address space, and refer to the corresponding location. The only real reason to have these in microcode is to avoid a timing error which can happen with some devices which have side effects when read. It is Lisp programs, not special microcode, which know the location and function of the registers in the keyboard, mouse, TV, and cable network interfaces. This makes the low-level I/O code just as flexible and easy to modify as the high level code.

There are also a couple of microcoded routines which speed up the drawing of characters in the TV memory. These do not do anything which could not be done in Lisp, but they are carefully hand-coded in microcode because we draw an awful lot of characters.

### High level:

Many programs perform simple stream-oriented (sequential characters) I/O. In order that these programs be kept device-independent, there is a standard definition of a "stream": a stream is a functional object which takes one required argument and one optional argument. The first argument is a symbol which is a "command" to the stream, such as "TYI", which means "input one character, and return it" and "TYO", which means "output one character". The character argument to the TYO command is passed in the second argument to the stream. There are several other standard optional stream operations, for several

purposes including higher efficiency. In addition particular devices can define additional operations for their own purposes.

Streams can be used for I/O to files in the file system, strings inside the Lisp Machine, the terminal, editor buffers, or anything else which is naturally represented as sequential characters.

For I/O which is of necessity device-dependent, such as the sophisticated operations performed on the TV by the editor, which include multiple blinkers and random access to the screen, special packages of Lisp functions are provided, and there is no attempt to be device-independent. (See documentation on the TV and network packages).

In general, we feel no regret at abandoning device independence in interactive programs which know they are using the display. The advantages to be gained from sophisticated high-bandwidth display-based interaction far outweigh the advantages of device-independence. This does mean that the Lisp machine is really not usable from other than its own terminal; in particular, it cannot be used remotely over the ARPANET.

## REPRESENTATION OF DATA:

A Lisp object in Maclisp or InterLisp is represented as an 18 bit pointer, and the datatype of the object is determined from the pointer; each page of memory can only contain objects of a single type. In the Lisp machine, Lisp objects are represented by a 5 bit datatype field, and a 24 bit pointer. (The Lisp machine virtual address space is 24 bits). There are a variety of datatypes (most of the 32 possible codes are now in use), which have symbolic names such as DTP-LIST, DTP-SYMBOL, DTP-FIXNUM, etc.

The Lisp machine data types are designed according to these criteria: There should be a wide variety of useful and flexible data types. Some effort should be made to increase the bit-efficiency of data representation, in order to improve performance. The programmer should be able to exercise control over the storage and representation of data, if he wishes. It must always be possible to take an anonymous piece of data and discover its type; this facilitates storage management. There should be as much type-checking and error-checking as feasible in the system.

Symbols are stored as four consecutive words, each of which contains one object. The words are termed the PRINT NAME cell, the VALUE cell, the FUNCTION

cell, and the PROPERTY LIST cell.  The PRINT NAME cell holds a string object, which is the printed representation of the symbol.  The PROPERTY LIST cell, of course, contains the property list, and the VALUE CELL contains the current value of the symbol (it is a shallow-binding system).  The FUNCTION cell replaces the task of the EXPR, SUBR, FEXPR, MACRO, etc.  properties in Maclisp.  When a form such as (FOO ARG1 ARG2) is evaluated, the object in FOO's function cell is applied to the arguments.  A symbol object has datatype DTP-SYMBOL, and the pointer is the address of these four words.

Storage of list structure is somewhat more complicated.  Normally a "list object" has datatype DTP-LIST, and the pointer is the address of a two word block;  the first word contains the CAR, and the second the CDR of the node.

However, note that since a Lisp object is only 29 bits (24 bits of pointer and 5 bits of data-type), there are three remaining bits in each word. Two of these bits are termed the CDR-CODE field, and are used to compress the storage requirement of list structure.  The four possible values of the CDR-CODE field are given the symbolic names CDR-NORMAL, CDR-ERROR, CDR-NEXT, and CDR-NIL. CDR-NORMAL indicates the two-word block described above.  CDR-NEXT and CDR-NIL are used to represent a list as a vector, taking only half as much storage as usual;  only the CARs are stored.  The CDR of each location is simply the next location, except for the last, whose CDR is NIL.  The primitive functions which create lists (LIST, APPEND, etc.)  create these compressed lists.  If RPLACD is done on such a list, it is automatically changed back to the conventional two-word representation, in a transparent way.

The idea is that in the first word of a list node the CAR is represented by 29 bits, and the CDR is represented by 2 bits.  It is a compressed pointer which can take on only 3 legal values:  to the symbol NIL, to the next location after the one it appears in, or indirect through the next location.  CDR-ERROR is used for words whose address should not ever be in a list object;  in a "full node", the first word is CDR-NORMAL, and the second is CDR-ERROR.  It is important to note that the cdr-code portion of a word is used in a different way from the data-type and pointer portion;  it is a property of the memory cell itself, not of the cell's contents.  A "list object" which is represented in compressed form still has data type DTP-LIST, but the cdr code of the word addressed by its pointer field is CDR-NEXT or CDR-NIL rather than CDR-NORMAL.

Number objects may have any of three datatypes.  "FIXNUMs", which are 24-bit signed integers, are represented by objects of datatype DTP-FIX, whose "pointer" parts are actually the value of the number.  Thus fixnums, unlike all other objects, do not require any "CONS"ed storage for their representation.

This speeds up arithmetic programs when the numbers they work with are reasonably small. Other types of numbers, such as floating point, BIGNUMs (integers of arbitrarily high precision), complex numbers, and so on, are represented by objects of datatype DTP-EXTENDED-NUMBER which point to a block of storage containing the details of the number. The microcode automatically converts between the different number representations as necessary, without the need for explicit declarations on the programmer's part.

There is also a datatype DTP-PDL-NUMBER, which is almost the same as DTP-EXTENDED-NUMBER. The difference is that pdl numbers can only exist in the pdl buffer (a memory internal to the machine which holds the most recent stack frames), and their blocks of storage are allocated in a special area. Whenever an object is stored into memory, if it is a pdl number its block of storage is copied, and an ordinary extended number is substituted. The idea of this is to prevent intermediate numeric results from using up storage and causing increased need for garbage collection. When the special pdl number area becomes full, all pdl numbers can quickly be found by scanning the pdl buffer. Once they have been copied out into ordinary numbers, the special area is guaranteed empty and can be reclaimed, with no need to garbage collect nor to look at other parts of memory. Note that these are not at all the same as pdl numbers in Maclisp; however, they both exist for the same reason.

The most important other data type is the array. Some problems are best attacked using data structures organized in the list-processing style of Lisp, and some are best attacked using the array-processing style of Fortran. The complete programming system needs both. As mentioned above, Lisp Machine arrays are augmented beyond traditional Lisp arrays in several ways. First of all, we have the ordinary arrays of Lisp objects, with one or more dimensions. Compact storage of positive integers, which may represent characters or other non-numeric entities, is afforded by arrays of 1-bit, 2-bit, 4-bit, 8-bit, or 16-bit elements.

For string-processing, there are string-arrays, which are usually one-dimensional and have 8-bit characters as elements. At the microcode level strings are treated the same as 8-bit arrays, however strings are treated differently by READ, PRINT, EVAL, and many other system and user functions. For example, they print out as a sequence of characters enclosed in quotes. The characters in a character string can be accessed and modified with the same array-referencing functions as one uses for any other type of array. Unlike arrays in other Lisp systems, Lisp machine arrays usually have only a single word

of overhead, so the character strings are quite storage-efficient.

There are a number of specialized types of arrays which are used to implement other data types, such as stack groups, internal system tables, and, most importantly, the refresh memory of the TV display as a two-dimensional array of bits.

An important additional feature of Lisp machine arrays is called "array leaders." A leader is a vector of Lisp objects, of user-specified size, which may be tacked on to an array. Leaders are a good place to remember miscellaneous extra information associated with an array. Many data structures consist of a combination of an array and a record (see below); the array contains a number of objects all of the same conceptual type, while the record contains miscellaneous items all of different conceptual types. By storing the record in the leader of the array, the single conceptual data structure is represented by a single actual object. Many data structures in Lisp-machine system programs work this way.

Another thing that leaders are used for is remembering the "current length" of a partially-populated array. By convention, array leader element number 0 is always used for this.

Many programs use data objects structured as "records"; that is, a compound object consisting of a fixed number of named sub-objects. To facilitate the use of records, the Lisp machine system includes a standard set of macros for defining, creating, and accessing record structures. The user can choose whether the actual representation is to be a Lisp list, an array, or an array-leader. Because this is done with macros, which translate record operations into the lower-level operations of basic Lisp, no other part of the system needs to know about records.

Since the reader and printer are written in Lisp and user-modifiable, this record-structure feature could easily be expanded into a full-fledged user-defined data type facility by modifying read and print to support input and output of record types.

Another data type is the "locative pointer." This is an actual pointer to a memory location, used by low-level system programs which need to deal with the guts of data representation. Taking CAR or CDR of a locative gets the contents of the pointed-to location, and RPLACA or RPLACD stores. It is possible to LAMBDA-bind the location. Because of the tagged architecture and highly-organized storage, it is possible to have a locative pointer into the middle of almost anything without causing trouble with the garbage collector.

## REPRESENTATION OF PROGRAMS:

In the Lisp Machine there are three representations for programs. Interpreted Lisp code is the slowest, but the easiest for programs to understand and modify. It can be used for functions which are being debugged, for functions which need to be understood by other functions, and for functions which are not worth the bother of compiling. A few functions, notably EVAL, will not work interpreted.

Compiled Lisp ("macrocode") is the main representation for programs. This consists of instructions in a somewhat conventional machine-language, whose unusual features will be described below. Unlike the case in many other Lisp systems, macrocode programs still have full checking for unbound variables, data type errors, wrong number of arguments to a function, and so forth, so it is not necessary to resort to interpreted code just to get extra checking to detect bugs. Often, after typing in a function to the editor, one skips the interpretation step and requests the editor to call the compiler on it, which only takes a few seconds since the compiler is always in the machine and only has to be paged in.

Compiled code on the Lisp Machine is stored inside objects called (for historical reasons) Function Entry Frames (FEFs). For each function compiled, one FEF is created, and an object of type DTP-FEF-POINTER is stored in the function cell of the symbol which is the name of the function. A FEF consists of some header information, a description of the arguments accepted by the function, pointers to external Lisp objects needed by the function (such as constants and special variables), and the macrocode which implements the function.

The third form of program representation is microcode. The system includes a good deal of hand-coded microcode which executes the macrocode instructions, implements the data types and the function-calling mechanism, maintains the paged virtual memory, does storage allocation and garbage collection, and performs similar systemic functions. The primitive operations on the basic data types, that is, CAR and CDR for lists, arithmetic for numbers, reference and store for arrays, etc. are implemented as microcode subroutines. In addition, a number of commonly-used Lisp functions, for instance GET and ASSQ, are hand-coded in microcode for speed.

In addition to this system-supplied microcode, there is a feature called micro compilation. Because of the simplicity and generality of the CONS microprocessor, it is feasible to write a compiler to compile user-written Lisp functions directly into microcode, eliminating the overhead of fetching and

interpreting macroinstructions. This can be used to boost performance by microcompiling the most critical routines of a program. Because it is done by a compiler rather than a system programmer, this performance improvement is available to everyone. The amount of speedup to be expected depends on the operations used by the program; simple low-level operations such as data transmission, byte extraction, integer arithmetic, and simple branching can expect to benefit the most. Function calling, and operations which already spend most of their time in microcode, such as ASSQ, will benefit the least. In the best case one can achieve a factor of about 20. In the worst case, maybe no speedup at all.

Since the amount of control memory is limited, only a small number of microcompiled functions can be loaded in at one time. This means that programs have to be characterized by spending most of their time in a small inner kernel of functions in order to benefit from microcompilation; this is probably true of most programs. There will be fairly hairy metering facilities for identifying such critical functions.

We do not yet have a microcompiler, but a prototype of one was written and heavily used as part of the Lisp machine simulator. It compiles for the PDP-10 rather than CONS, but uses similar techniques and a similar interface to the built-in microcode.

In all three forms of program, the flexibility of function calling is augmented with generalized LAMBDA-lists. In order to provide a more general and flexible scheme to replace EXPRs vs. FEXPRs vs. LEXPRs, a syntax borrowed from Muddle and Conniver is used in LAMBDA lists. In the general case, there are an arbitrary number of REQUIRED parameters, followed by an arbitrary number of OPTIONAL parameters, possibly followed by one REST parameter. When a function is APPLIED to its arguments, first of all the required formal parameters are paired off with arguments; if there are fewer arguments than required parameters, an error condition is caused. Then, any remaining arguments are paired off with the optional parameters; if there are more optional parameters than arguments remaining, then the rest of the optional parameters are initialized in a user-specified manner. The REST parameter is bound to a list, possibly NIL, of all arguments remaining after all OPTIONAL parameters are bound. To avoid CONSing, this list is actually stored on the pdl; this means that you have to be careful how you use it, unfortunately. It is also possible to control which arguments are evaluated and which are quoted.

Normally, such a complicated calling sequence would entail an

unacceptable amount of overhead.  Because this is all implemented by microcode, and because the simple, common cases are special-cased, we can provide these advanced features and still retain the efficiency needed in a practical system.

We will now discuss some of the issues in the design of the macrocode instruction set.  Each macroinstruction is 16 bits long;  they are stored two per word.  The instructions work in a stack-oriented machine.  The stack is formatted into frames;  each frame contains a bunch of arguments, a bunch of local variable value slots, a push-down stack for intermediate results, and a header which gives the function which owns the frame, links this frame to previous frames, remembers the program counter and flags when this frame is not executing, and may contain "additional information" used for certain esoteric purposes.  Originally this was intended to be a spaghetti stack, but the invention of closures and stack-groups (see the control-structure section), combined with the extreme complexity of spaghetti stacks, made us decide to use a simple linear stack.  The current frame is always held in the pdl buffer, so accesses to arguments and local variables do not require memory references, and do not have to make checks related to the garbage collector, which improves performance.  Usually several other frames will also be in the pdl buffer.

The macro instruction set is bit-compact.  The stack organization and Lisp's division of programs into small, separate functions means that address fields can be small.  The use of tagged data types, powerful generic operations, and easily-called microcoded functions makes a single 16-bit macro instruction do the work of several instructions on a conventional machine such as a PDP-10.

The primitive operations which are the instructions which the compiler generates are higher-level than the instructions of a conventional machine.  They all do data type checks;  this provides more run-time error checking than in Maclisp, which increases reliability.  But it also eliminates much of the need to make declarations in order to get efficient code.  Since a data type check is being made, the "primitive" operations can dynamically decide which specific routine is to be called.  This means that they are all "generic", that is, they work for all data types where they make sense.

The operations which are regarded as most important, and hence are easiest for macrocode to do, are data transmission, function calling, conditional testing, and simple operations on primitive types, that is, CAR, CDR, CADR, CDDR, RPLACA, and RPLACD, plus the usual arithmetic operations and comparisons.  More complex operations are generally done by "miscellaneous" instructions, which call microcoded subroutines, passing arguments on the temporary-results stack.

There are three main kinds of addressing in macrocode.  First, there is

implicit addressing of the top of the stack. This is the usual way that operands get from one instruction to the next.

Second, there is the source field (this is sometimes used to store results, but I will call it a source anyway). The source can address any of the following: Up to 64 arguments to the current function. Up to 64 local variables of the current function. The last result, popped off the stack. One of several commonly-used constants (e.g. NIL) stored in a system-wide constants area. A constant stored in the FEF of this function. A value cell or a function cell of a symbol, referenced by means of an invisible pointer in the FEF; this mode is used to reference special variables and to call other functions.

Third, there is the destination field, which specifies what to do with the result of the instruction. The possibilities are: Ignore it, except set the indicators used by conditional branches. Push it on the stack. Pass it as an argument. Return it as the value of this function. Cons up a list.

There are five types of macroinstructions, which will be described. First, there are the data transmission instructions, which take the source and MOVE it to the destination, optionally taking CAR, CDR, CAAR, CADR, CDAR, or CDDR in the process. Because of the powerful operations that can be specified in the destination, these instructions also serve as argument-passing, function-exiting, and list-making instructions.

Next we have the function calling instructions. The simpler of the two is CALL0, call with no arguments. It calls the function indicated by its source, and when that function returns, the result is stored in the destination. The microcode takes care of identifying what type of function is being called, invoking it in the appropriate way, and saving the state of the current function. It traps to the interpreter if the called function is not compiled.

The more complex function call occurs when there are arguments to be passed. The way it works is as follows. First, a CALL instruction is executed. The source operand is the function to be called. The beginnings of a new stack frame are constructed at the end of the current frame, and the function to be called is remembered. The destination of the CALL instruction specifies where the result of the function will be placed, and it is saved for later use when the function returns. Next, instructions are executed to compute the arguments and store them into the destination NEXT-ARGUMENT. This causes them to be added to the new stack frame. When the last argument is computed, it is stored into the destination LAST-ARGUMENT, which stores it in the new stack frame and then activates the call. The function to be called is analyzed, and the arguments are bound to the formal parameters (usually the arguments are already in the correct

slots of the new stack frame). Because the computation of the arguments is introduced by a CALL instruction, it is easy to find out where the arguments are and how many there are. The new stack frame becomes current and that function begins execution. When it returns, the saved destination of the CALL instruction is retrieved and the result is stored. Note that by using a destination of NEXT-ARGUMENT or LAST-ARGUMENT function calls may be nested. By using a destination of RETURN the result of one function may become the result of its caller.

The third class of macro instructions consists of a number of common operations on primitive data types. These instructions do not have an explicit destination, in order to save bits, but implicitly push their result (if any) onto the stack. This sometimes necessitates the generation of an extra MOVE instruction to put the result where it was really wanted. These instructions include: Operations to store results from the pdl into the "source". The basic arithmetic and bitwise boolean operations. Comparison operations, including EQ and arithmetic comparison, which set the indicators for use by conditional branches. Instructions which set the "source" operand to NIL or zero. Iteration instructions which change the "source" operand using CDR, CDDR, 1+, or 1- (add or subtract one). Binding instructions which lambda-bind the "source" operand, then optionally set it to NIL or to a value popped off the stack. And, finally, an instruction to push its effective address on the stack, as a locative pointer.

The fourth class of macro instructions are the branches, which serve mainly for compiling COND. Branches contain a self-relative address which is transferred to if a specified condition is satisfied. There are two indicators, which tell if the last result was NIL, and if it was an atom, and the state of these indicators can be branched on; there is also an unconditional branch, of course. For branches more than 256 half-words away, there is a double-length long-branch instruction. An interesting fact is that there are not really any indicators; it turns out to be faster just to save the last result in its entirety, and compare it against NIL or whatever when that is needed by a branch instruction. It only has to be saved from one instruction to the immediately following one.

The fifth class of macro instructions is the "miscellaneous function." This selects one of 512 microcoded functions to be called, with arguments taken from results previously pushed on the stack. A destination is specified to receive the result of the function. In addition to commonly-used functions such as GET, CONS, CDDDDR, REMAINDER, and ASSQ, miscellaneous functions include sub-primitives (discussed above), and instructions which are not as commonly used

as the first four classes, including operations such as array-accessing, consing
up lists, un-lambda-binding, special funny types of function calling, etc.

The way consing-up of lists works is that one first does a miscellaneous
function saying "make a list N long". One then executes N instructions with
destination NEXT-LIST to supply the elements of the list. After the Nth such
instruction, the list-object magically appears on the top of the stack. This
saves having to make a call to the function LIST with a variable number of
arguments.

Another type of "instruction set" used with macrocode is the Argument
Description List, which is executed by a different microcoded interpreter at the
time a function is entered. The ADL contains one entry for each argument which
the function expects to be passed, and for each auxiliary variable. It contains
all relevant information about the argument: whether it is required, optional,
or rest, how to initialize it if it is not provided, whether it is local or
special, datatype checking information, and so on. Sometimes the ADL can be
dispensed with if the "fast argument option" can be used instead; this helps
save time and memory for small, simple functions. The fast-argument option is
used when the optional arguments and local variables are all to be initialized to
NIL, there are not too many of them, there is no data-type checking, and the
usage of special variables is not too complicated. The selection of the
fast-argument option, if appropriate, is automatically made by the system, so the
user need not be concerned with it. The details can be found in the FORMAT
document.


## CONTROL STRUCTURES:

Function calling. Function calling is, of course, the basic main control
structure in Lisp. As mentioned above, Lisp machine function calling is made
fast through the use of microcode and augmented with optional arguments, rest
arguments, multiple return values, and optional type-checking of arguments.

CATCH and THROW. CATCH and THROW are a Maclisp control structure which
will be mentioned here since they may be new to some people. CATCH is a way of
marking a particular point in the stack of recursive function invocations. THROW
causes control to be unwound to the matching CATCH, automatically returning
through the intervening function calls. They are used mainly for handling errors
and unusual conditions. They are also useful for getting out of a hairy piece of

code when it has discovered what value it wants to return;  this applies
particularly to nested loops.

Closures.  The LISP machine contains a data-type called "closure" which
is used to implement "full funarging".  By turning a function into a closure, it
becomes possible to pass it as an argument with no worry about naming conflicts,
and to return it as a value with exactly the minimum necessary amount of binding
environment being retained, solving the classical "funarg problem".  Closures are
implemented in such a way that when they are not used the highly speed- and
storage-efficient shallow binding variable scheme operates at full efficiency,
and when they are used things are slowed down only slightly.  The way one creates
a closure is with a form such as:


        (CLOSURE '(FOO-PARAM FOO-STATE)
                 (FUNCTION FOO-BAR))


The function could also be written directly in place as a
LAMBDA-expression, instead of referring to the externally defined FOO-BAR.  The
variables FOO-PARAM and FOO-STATE are those variables which are used free by
FOO-BAR and are intended to be "closed".  That is, these are the variables whose
binding environment is to be fixed to that in effect at the time the closure is
created.  The explicit declaration of which variables are to be closed allows the
implementation to have high efficiency, since it does not need to save the whole
variable-binding environment, almost all of which is useless.  It also allows the
programmer to explicitly choose for each variable whether it is to be dynamically
bound (at the point of call) or statically bound (at the point of creation of the
closure), a choice which is not conveniently available in other languages.  In
addition the program is clearer because the intended effect of the closure is
made manifest by listing the variables to be affected.


Here is an example, in which the closure feature is used to solve a
problem presented in "LAMBDA - The Ultimate Imperative" [LAMBDA].  The problem is
to write a function called GENERATE-SQRT-OF-GIVEN-EXTRA-TOLERANCE, which is to
take one argument, which is the factor by which the tolerance is to be increased,
and return a function which takes square roots with that much more tolerance than
usual, whatever "usual" is later defined to be.  You are given a function SQRT
which makes a free reference to EPSILON, which is the tolerance it demands of the
trial solution.  The reason this example presents difficulties to various
languages is that the variable EPSILON must be bound at the point of call (i.e.

dynamically scoped), while the variable FACTOR must be bound at the point of
creation of the function (i.e. lexically scoped). Thus the programmer must have
explicit control over how the variables are bound.

```
(DEFUN GENERATE-SQRT-OF-GIVEN-EXTRA-TOLERANCE (FACTOR)
    (CLOSURE '(FACTOR)
            (FUNCTION
                (LAMBDA (X)
                    ((LAMBDA (EPSILON) (SQRT X))
                     (* EPSILON FACTOR))))))
```

The function, when called, rebinds EPSILON to FACTOR times its current value,
then calls SQRT. The value of FACTOR used is that in effect when the closure was
created, i.e. the argument to GENERATE-SQRT-OF-GIVEN-EXTRA-TOLERANCE.

The way closures are implemented is as follows. For each variable to be
closed an "external value cell" is created, which is a CONSed up free-storage
cell which contains the variable's value when it is at that level of binding.
Because this cell is CONSed up, it can be retained as long as necessary, just
like any other data, and unlike cells in a stack. Because it is a cell, if the
variable is SETQed the new value is seen by all the closures that should see it.
The association between the symbol which is the name of the variable and this
value cell is of the shallow-binding type, for efficiency; an invisible pointer
(see the storage organization section) in the normal (internal) value cell
supplies the connection, eliminating the overhead of searching stack frames or
a-lists. If at the time the closure is created an external value cell already
exists for a variable, that one is used instead of creating a new one. Thus all
closures at the same "level of binding" use the same value cell, which is the
desired semantics.

The CLOSURE function returns an object of type DTP-CLOSURE, which
contains the function to be called and, for each variable closed over, locative
pointers to its internal and external value cells.

When a closure is invoked as a function, the variables mentioned in the
closure are bound to invisible pointers to their external value cells; this puts
these variables into the proper binding environment. The function contained in
the closure is then invoked in the normal way. When the variables happen to be
referred to, the invisible pointers are automatically followed to the external
value cells. If one of the closed variables is then bound by some other

function, the external value cell pointer is saved away on the binding stack, like any saved variable value, and the variable reverts to normal nonclosed status. When the closed function returns, the bindings of the closed variables are restored just like any other variables bound by the function.

Note the economy of mechanism. Almost all of the system is completely unaffected by and unaware of the existence of closures; the invisible pointer mechanism takes care of things. The retainable binding environments are allocated through the standard CONS operation. The switching of variables between normal and "closed" status is done through the standard binding operation. The operations used by a closed function to access the closed variables are the same as those used to access ordinary variables; closures are called in the same way as ordinary functions. Closures work just as well in the interpreter as in the compiler. An important thing to note is the minimality of CONSing in closures. When a closure is created, some CONSing is done; external value cells and the closure-object itself must be created, but there is no extra "overhead". When a closure is called, no CONSing happens.

One thing to note is that in the compiler closed variables have to be declared "special". This is a general feature of the Maclisp and Lisp machine compilers, that by default variables are local, which means that they are lexically bound, only available to the function in which they are bound, and implemented not with atomic symbols, but simply as slots in the stack. Variables that are declared special are implemented with shallow-bound atomic symbols, identical to variables in the interpreter, and have available either dynamic binding or closure binding. They are somewhat less efficient since it takes two memory references to access them and several to bind them.

Stack groups. The stack group is a type of Lisp object useful for implementation of certain advanced control structures such as coroutines, asynchronous processes, and generators. A stack group is similar to a process (or fork or job or task or control-point) in a time-sharing system; it contains such state information as the "regular" and "special" (binding) PDLs and various internal registers. At all times there is one stack group running on the machine.

Control may be passed between stack groups in several ways (not all of which exist yet on our prototype machine). A stack-group may be called like a function; when it wants to return it can do a %STACK-GROUP-RETURN which is different from an ordinary function return in that the state of the stack group remains unchanged; the next time it is called it picks up from where it left

off.  This is good for generator-like applications;  each time
%STACK-GROUP-RETURN is done, a value is emitted from the generator, and as a
side-effect execution is suspended until the next time the generator is called.
%STACK-GROUP-RETURN is analogous to the ADIEU construct in CONNIVER.

Control can simply be passed explicitly from one stack group to another,
coroutine-style.  Alternatively, there can be a scheduler stack-group which
invokes other stack groups when their requested scheduling conditions are
satisfied.

Interrupts cause control of the machine to be transferred to an
interrupt-handler stack group.  Essentially this is a forced stack group call
like those calls described above.  Similarly, when the microcode detects an error
the current stack group is suspended and control is passed to an error-handling
stack group.  The state of the stack group that got the error is left exactly as
it was when the error occurred, undisturbed by any error-handling operations.
This facilitates error analysis and recovery.

When the machine is started, an "initial" stack group becomes the current
stack group, and is forced to call the first function of Lisp.

Note that the same scheduler-driven stack-group switching mechanism can
be used both for user programs which want to do parallel computations, and for
system programming purposes such as the handling of network servers and
peripheral handlers.

Each stack group has a call-state and a calling-stack-group variable,
which are used in maintaining the relations between stack groups.  A stack group
also has some option flags controlling whether the system tries to keep different
stack groups' binding environments distinct by undoing the special variable
bindings of the stack group being left and redoing the bindings of the stack
group being entered.

Stack groups are created with the function MAKE-STACK-GROUP, which takes
one main argument, the "name" of the stack group.  This is used only for
debugging, and can be any mnemonic symbol.  It returns the stack group, i.e., a
Lisp object with data type DTP-STACK-GROUP.  Optionally the sizes of the pdls may
be specified.

The function STACK-GROUP-PRESET is used to initialize the state of a
stack group:  the first argument is the stack group, the second is a function to
be called when the stack group is invoked, and the rest are arguments to that
function.  Both PDLs are made empty.  The stack group is set to the
AWAITING-INITIAL-CALL state.  When it is activated, the specified function will
find that it has been called with the specified arguments.  If it should return

in the normal way, (i.e. the stack group "returns off the top", the stack group
will enter a "used up" state and control will revert to the calling stack group.
Normally, the specified function will use %STACK-GROUP-RETURN several times;
otherwise it might as well have been called directly rather than in a stack
group.

One important difference between stack groups and other means proposed to
implement similar features is that the stack group scheme involves no loss of
efficiency in normal computation.  In fact, the compiler, the interpreter, and
even the runtime function-calling mechanism are completely unaware of the
existance of stack groups.

## STORAGE ORGANIZATION:

Incremental Garbage Collection.  The Lisp machine will use a real-time,
incremental, compacting garbage collector.  Real-time means that CONS (or related
functions) never delay Lisp execution for more than a small, bounded amount of
time.

This is very important in a machine with a large address space, where a
traditional garbage collection could bring everything to a halt for several
minutes.  The garbage collector is incremental, i.e. garbage collection is
interleaved with execution of the user's program;  every time you call CONS the
garbage collection proceeds for a few steps.  Copying can also be triggered by a
memory reference which fetches a pointer to data which has not yet been copied.
The garbage collector compactifies in order to improve the paging
characteristics.

The basic algorithm is described in a paper by Henry Baker [GC].  We have
not implemented it yet, but design is proceeding and most of the necessary
changes to the microcode have already been made.  It is much simpler than
previous methods of incremental garbage collection in that only one process is
needed;  this avoids interlocking and synchronization problems, which are often
very difficult to debug.

Areas.  Storage in the Lisp machine is divided into "areas."  Each area
contains related objects, of any type.  Since unlike PDP-10 Lisps we do not
encode the data type in the address, we are free to use the address to encode the
area.  Areas are intended to give the user control over the paging behavior of

his program, among other things. By putting related data together, locality can be greatly increased. Whenever a new object is created, for instance with CONS, the area to be used can optionally be specified. There is a default Working Storage area which collects those objects which the user has not chosen to control explicitly.

Areas also give the user a handle to control the garbage collector. Some areas can be declared to be "static", which means that they change slowly and the garbage collector should not attempt to reclaim any space in them. This can eliminate a lot of useless copying. All pointers out of a static area can be collected into an "exit vector", eliminating any need for the garbage collector to look at that area. As an important example, an English-language dictionary can be kept inside the Lisp without adversely affecting the speed of garbage collection. A "static" area can be explicitly garbage-collected at infrequent intervals when it is believed that that might be worthwhile.

Each area can potentially have a different storage discipline, a different paging algorithm, and even a different data representation. The microcode will dispatch on an attribute of the area at the appropriate times. The structure of the machine makes the performance cost of these features negligible; information about areas is stored in extra bits in the memory mapping hardware where it can be quickly dispatched on by the microcode. These dispatches usually have to be done anyway to make the garbage collector work, and to implement invisible pointers.

Invisible Pointers. An invisible pointer is similar to an indirect address word on a conventional computer except the indirection is specified in the data instead of in the instruction. A reference to a memory location containing an invisible pointer is automatically altered to use the location pointed to by the invisible pointer. The term "invisible" refers to the fact that the presence of such pointers is not visible to most of the system, since they are handled by the lowest-level memory-referencing operations. The invisible pointer feature does not slow anything down too much, because it is part of the data type checking that is done anyway (this is one of the benefits of a tagged architecture). A number of advanced features of the Lisp machine depend upon invisible pointers for their efficient implementation.

Closures use invisible pointers to connect internal value cells to external value cells. This allows the variable binding scheme to be altered from normal shallow binding to allocated-value-cell shallow binding when closures are being used, without altering the normal operation of the machine when closures

are not being used.  At the same time the slow-down when closures are used amounts to only 2 microseconds per closed-variable reference, the time needed to detect and follow the invisible pointer.

Invisible pointers are necessary to the operation of the cdr-coded compressed list scheme.  If an RPLACD is done to a compressed list, the list can no longer be represented in the compressed form.  It is necessary to allocate a full 2-word cons node and use that in its place.  But, it is also necessary to preserve the identity (with respect to EQ) of the list.  This is done by storing an invisible pointer in the original location of the compressed list, pointing to the uncompressed copy.  Then the list is still represented by its original location, preserving EQ-ness, but the CAR and CDR operations follow the invisible pointer to the new location and find the proper car and cdr.

This is a special case of the more general use of invisible pointers for "forwarding" references from an old representation of an object to a new one. For instance, there is a function to increase the size of an array.  If it cannot do it in place, it makes a new copy and leaves behind an invisible pointer.

The exit-vector feature uses invisible pointers.  One may set up an area to have the property that all references from inside that area to objects in other areas are collected into a single exit-vector.  A location which would normally contain such a reference instead contains an invisible pointer to the appropriate slot in the exit vector.  Operations on this area all work as before, except for a slight slow-down caused by the invisible pointer following.  It is also desirable to have automatic checking to prevent the creation of new outside references;  when an attempt is made to store an outside object into this area execution can trap to a routine which creates a new exit vector entry if necessary and stores an invisible pointer instead.  The reason for exit vectors is to speed up garbage collection by eliminating the need to swap in all of the pages of the area in order to find and relocate all its references to outside objects.

The macrocode instruction set relies on invisible pointers in order to access the value cells of "special" (non-local) variables and the function cells of functions to be called.

Certain system variables stored in the microcode scratchpad memory are made available to Lisp programs by linking the value cells of appropriately-named Lisp symbols to the scratchpad memory locations with invisible pointers.  This makes it possible not only to read and write these variables, but also to lambda-bind them.  In a similar fashion, invisible pointers could be used to link two symbols' value cells together, in the fashion of MicroPlanner but with much

greater efficiency.


## THE EDITOR:

The Lisp machine system includes an advanced real-time display oriented editor, which is written completely in Lisp. The design of this editor drew heavily on our experience with the EMACS editor (and its predecessors) on the PDP-10. The high-speed display and fast response time of the Lisp machine are crucial to the success of the editor.

The TV display is used to show a section of the text buffer currently being edited. When the user types a normal printing character on the keyboard, that character is inserted into his buffer, and the display of the buffer is updated; you see the text as you type it in. When using an editor, most of the user's time is spent in typing in text; therefore, this is made as easy as possible. Editing operations other than the insertion of single characters are invoked by control-keys, i.e. by depressing the CONTROL and/or META shift keys, along with a single character. For example, the command to move the current location for typein in the buffer (the "point") backward is Control-B (B is mnemonic for Backward); the command to move to the next line is Control-N. There are many more advanced commands, which know how to interpret the text as words or as the printed representation of Lisp data structure; Meta-F moves forward over an English word, and Control-Meta-F moves forward over a Lisp expression (an atom or a list).

The real-time display-oriented type of editor is much easier to use than traditional text editors, because you can always see exactly what you are doing. A new user can sit right down and type in text. However, this does not mean that there can be no sophisticated commands and macros. Very powerful operations are provided in the Lisp machine editor. Self-documentation features exist to allow the user to ask what a particular key does before trying it, and to ask what keys contain a given word in their description. Users can write additional commands, in Lisp, and add them to the editor's command tables.

The editor knows how much a line should be indented in a Lisp program in order to reflect the level of syntactic nesting. When typing in Lisp code, one uses the Linefeed key after typing in a line to move to the next line and automatically indent it by the right amount. This serves the additional purpose of instantly pointing out errors in numbers of parentheses.

The editor can be used as a front end to the Lisp top level loop. This provides what can be thought of as very sophisticated rubout processing. When the user is satisfied that the form as typed is correct, he can activate it, allowing Lisp to read in the form and evaluate it. When Lisp prints out the result, it is inserted into the buffer at the right place. Simple commands are available to fetch earlier inputs, for possible editing and reactivation.

In addition to commands from the keyboard, the mouse can be used to point to parts of the buffer, and to give simple editing commands. The use of mice for text editing was originated at SRI, and has been refined and extended at XEROX-PARC.

The character-string representation of each function in a program being worked on is stored in its own editor buffer. One normally modifies functions by editing the character-string form, then typing a single-character command to read it into Lisp, replacing the old function. Compilation can optionally be requested. The advantage of operating on the character form, rather than directly on the list structure, is that comments and the user's chosen formatting of the code are preserved; in addition, the editor is easier to use because it operates on what you see on the display. There are commands to store sets of buffers into files, and to get them back again.

The editor has the capability to edit and display text in multiple fonts, and many other features too numerous to mention here.

## CURRENT STATUS (August 1977)

The original prototype CONS machine was designed and built somewhat more than two years ago. It had no memory and no I/O capability, and remained pretty much on the back burner while software was developed with a simulator on the PDP-10 (the simulator executed the Lisp machine macro instruction set, a function now performed by CONS microcode.) Microprogramming got under way a little over a year ago, and in the beginning of 1977 the machine got memory, a disk, and a terminal.

We now have an almost-complete system running on the prototype machine. The major remaining "holes" are the lack of a garbage collector and the presence of only the most primitive error handling. Also, floating-point and big-integer numbers and microcompilation have been put off until the next machine. The system includes almost all the functions of Maclisp, and quite a few new ones.

The machine is able to page off of its disk, accept input from the keyboard and the mouse, display on the TV, and do I/O to files on the PDP-10. The display editor is completely working, and the compiler runs on the machine, so the system is quite usable for typing in, editing, compiling, and debugging Lisp functions.

As a demonstration of the system, and a test of its capabilities, two large programs have been brought over from the PDP-10. William Woods's LUNAR English-language data-base query system was converted from InterLisp to Maclisp, thence to Lisp machine Lisp. On the Lisp machine it runs approximately 3 times as fast as in Maclisp on the KA-10, which in turn is 2 to 4 times as fast as in InterLisp. Note that the Lisp machine time is elapsed real time, while the PDP-10 times are virtual run times as given by the operating system and do not include the delays due to timesharing.

Most of the Macsyma symbolic algebraic system has been converted to the Lisp machine; nearly all the source files were simply compiled without any modifications. Most of Macsyma works except for some things that require bignums. The preliminary speed is the same as on the KA-10, but a number of things have not been optimally converted. (This speed measurement is, again, elapsed time on the Lisp machine version versus reported run time on the KA-10 time sharing system. Thus, paging and scheduling overhead in the KA-10 case are not counted in this measurement.)

LUNAR (including the dictionary) and Macsyma can reside together in the Lisp machine with plenty of room left over; either program alone will not entirely fit in a PDP-10 address space.

The CONS machine is currently being redesigned, and a new machine will be built soon, replacing our present prototype. The new machine will have larger sizes for certain internal memories, will incorporate newer technology, will have greatly improved packaging, and will be faster. It will fit entirely in one cabinet and will be designed for ease of construction and servicing. In late 1977 and early 1978 we plan to build seven additional machines and install them at the MIT AI Lab. During the fall of 1977 we plan to finish the software, bringing it to a point where users can be put on the system. User experience with the Lisp machine during 1978 should result in improvement and cleaning up of the software and documentation, and should give us a good idea of the real performance to be expected from the machine. At that time we will be able to start thinking about ways to make Lisp machines available to the outside world.

## REFERENCES:

CONS:   Steele, Guy L. "Cons", not yet published.  This is a revision of
             Working paper 80, CONS by Tom Knight

GC:   Baker, Henry, "List Processing in Real Time on a Serial Computer",
             Working Paper 139

LAMBDA:   Steele, Guy L. "LAMBDA - The Ultimate Imperative", Artificial
Intelligence
             Memo 353

MOUSE:   See extensive publications by Englebart and group at SRI.

```
 1  /*---------------------------------------------------------------------
 2    key.c :          Key Management Engine for BSD
 3
 4    Copyright 1995 by Bao Phan,  Randall Atkinson, & Dan McDonald,
 5    All Rights Reserved.  All Rights have been assigned to the US
 6    Naval Research Laboratory (NRL).  The NRL Copyright Notice and
 7    License governs distribution and use of this software.
 8
 9    Patents are pending on this technology.  NRL grants a license
10    to use this technology at no cost under the terms below with
11    the additional requirement that software, hardware, and
12    documentation relating to use of this technology must include
13    the note that:
14         This product includes technology developed at and
15      licensed from the Information Technology Division,
16      US Naval Research Laboratory.
17
18    ----------------------------------------------------------------*/
19  /*---------------------------------------------------------------------
20  #   @(#)COPYRIGHT   1.1a (NRL) 17 August 1995
21
22  COPYRIGHT NOTICE
23
24  All of the documentation and software included in this software
25  distribution from the US Naval Research Laboratory (NRL) are
26  copyrighted by their respective developers.
27
28  This software and documentation were developed at NRL by various
29  people.  Those developers have each copyrighted the portions that they
30  developed at NRL and have assigned All Rights for those portions to
31  NRL.  Outside the USA, NRL also has copyright on the software
32  developed at NRL. The affected files all contain specific copyright
33  notices and those notices must be retained in any derived work.
34
35  NRL LICENSE
36
37  NRL grants permission for redistribution and use in source and binary
38  forms, with or without modification, of the software and documentation
39  created at NRL provided that the following conditions are met:
40
41  1. Redistributions of source code must retain the above copyright
42     notice, this list of conditions and the following disclaimer.
43  2. Redistributions in binary form must reproduce the above copyright
44     notice, this list of conditions and the following disclaimer in the
45     documentation and/or other materials provided with the distribution.
46  3. All advertising materials mentioning features or use of this software
47     must display the following acknowledgement:
48
49      This product includes software developed at the Information
50      Technology Division, US Naval Research Laboratory.
51
52  4. Neither the name of the NRL nor the names of its contributors
53     may be used to endorse or promote products derived from this software
54     without specific prior written permission.
55
56  THE SOFTWARE PROVIDED BY NRL IS PROVIDED BY NRL AND CONTRIBUTORS ``AS
57  IS'' AND ANY EXPRESS OR IMPLIED WARRANTIES, INCLUDING, BUT NOT LIMITED
58  TO, THE IMPLIED WARRANTIES OF MERCHANTABILITY AND FITNESS FOR A
59  PARTICULAR PURPOSE ARE DISCLAIMED.  IN NO EVENT SHALL NRL OR
60  CONTRIBUTORS BE LIABLE FOR ANY DIRECT, INDIRECT, INCIDENTAL, SPECIAL,
61  EXEMPLARY, OR CONSEQUENTIAL DAMAGES (INCLUDING, BUT NOT LIMITED TO,
62  PROCUREMENT OF SUBSTITUTE GOODS OR SERVICES; LOSS OF USE, DATA, OR
63  PROFITS; OR BUSINESS INTERRUPTION) HOWEVER CAUSED AND ON ANY THEORY OF
64  LIABILITY, WHETHER IN CONTRACT, STRICT LIABILITY, OR TORT (INCLUDING
65  NEGLIGENCE OR OTHERWISE) ARISING IN ANY WAY OUT OF THE USE OF THIS
66  SOFTWARE, EVEN IF ADVISED OF THE POSSIBILITY OF SUCH DAMAGE.
67
```

1

```
 68   The views and conclusions contained in the software and documentation
 69   are those of the authors and should not be interpreted as representing
 70   official policies, either expressed or implied, of the US Naval
 71   Research Laboratory (NRL).
 72
 73   ------------------------------------------------------------------*/
 74
 75   #include <sys/types.h>
 76   #include <sys/param.h>
 77   #include <sys/proc.h>
 78   #include <sys/mbuf.h>
 79   #include <sys/socket.h>
 80   #include <sys/socketvar.h>
 81   #include <sys/time.h>
 82   #include <sys/kernel.h>
 83   #include <net/raw_cb.h>
 84   #include <net/if.h>
 85   #include <net/if_types.h>
 86   #include <net/if_dl.h>
 87   #include <net/route.h>
 88   #include <netinet/in.h>
 89   #include <netinet/in_var.h>
 90   #include <netinet/if_ether.h>
 91
 92   #include <netinet6/in6.h>
 93   #include <netinet6/in6_var.h>
 94   #include <netinet6/ipsec.h>
 95   #include <netinet6/key.h>
 96   #include <netinet6/in6_debug.h>
 97
 98   #define MAXHASHKEYLEN (2 * sizeof(int) + 2 * sizeof(struct
      sockaddr_in6))
 99
100   /*
101    *  Not clear whether these values should be
102    *  tweakable at kernel config time.
103    */
104   #define KEYTBLSIZE 61
105   #define KEYALLOCTBLSIZE 61
106   #define SO2SPITBLSIZE 61
107
108   /*
109    *  These values should be tweakable...
110    *  perhaps by using sysctl
111    */
112
113   #define MAXLARVALTIME 240;     /* Lifetime of a larval key table entry */
114   #define MAXKEYACQUIRE 1;       /* Max number of key acquire messages sent
      */
115                                  /*   per destination address
                                     */
116   #define MAXACQUIRETIME 15;     /* Lifetime of acquire message */
117
118   /*
119    *  Key engine tables and global variables
120    */
121
122   struct key_tblnode keytable[KEYTBLSIZE];
123   struct key_allocnode keyalloctbl[KEYALLOCTBLSIZE];
124   struct key_so2spinode so2spitbl[SO2SPITBLSIZE];
125
126   struct keyso_cb keyso_cb;
127   struct key_tblnode nullkeynode   = { 0, 0, 0, 0, 0 };
128   struct key_registry *keyregtable;
129   struct key_acquirelist *key_acquirelist;
130   u_long maxlarvallifetime = MAXLARVALTIME;
131   int maxkeyacquire = MAXKEYACQUIRE;
```

2

```
132   u_long maxacquiretime = MAXACQUIRETIME;
133
134   extern void dump_secassoc();
135
136
137   /*-----------------------------------------------------------------
138    * (temporary) Dump a data buffer
139    -----------------------------------------------------------------*
      /
140
141   void
142   dump_buf(buf, len)
143       char *buf;
144       int len;
145   {
146     int i;
147
148     printf("buf=0x%x len=%d:\n", buf, len);
149     for (i = 0; i < len; i++) {
150       printf("0x%x ", (u_int8)*(buf+i));
151     }
152     printf("\n");
153   }
154
155
156   /*-----------------------------------------------------------------
157    * (temporary) Dump a key tblnode structrue
158    -----------------------------------------------------------------*
      /
159
160   void
161   dump_keytblnode(ktblnode)
162       struct key_tblnode *ktblnode;
163   {
164     if (!ktblnode) {
165       printf("NULL key table node pointer!\n");
166       return;
167     }
168     printf("solist=0x%x ", ktblnode->solist);
169     printf("secassoc=0x%x ", ktblnode->secassoc);
170     printf("next=0x%x\n", ktblnode->next);
171   }
172
173
174   /*-----------------------------------------------------------------
175    * key_secassoc2msghdr():
176    *       Copy info from a security association into a key message buffer.
177    *       Assume message buffer is sufficiently large to hold all security
178    *       association information including src, dst, from, key and iv.
179    -----------------------------------------------------------------*
      /
180   int
181   key_secassoc2msghdr(secassoc, km, keyinfo)
182       struct ipsec_assoc *secassoc;
183       struct key_msghdr *km;
184       struct key_msgdata *keyinfo;
185   {
186     char *cp;
187     DPRINTF(IDL_GROSS_EVENT, ("Entering key_secassoc2msghdr\n"));
188
189     if ((km == 0) || (keyinfo == 0) || (secassoc == 0))
190       return(-1);
191
192     km->type = secassoc->type;
193     km->state = secassoc->state;
194     km->label = secassoc->label;
195     km->spi = secassoc->spi;
```

3

```
196     km->keylen = secassoc->keylen;
197     km->ivlen = secassoc->ivlen;
198     km->algorithm = secassoc->algorithm;
199     km->lifetype = secassoc->lifetype;
200     km->lifetime1 = secassoc->lifetime1;
201     km->lifetime2 = secassoc->lifetime2;
202
203     /*
204      *  Stuff src/dst/from/key/iv in buffer after
205      *  the message header.
206      */
207     cp = (char *)(km + 1);
208
209 #define ROUNDUP(a) \
210     ((a) > 0 ? (1 + (((a) - 1) | (sizeof(long) - 1))) : sizeof(long))
211 #define ADVANCE(x, n) \
212     { x += ROUNDUP(n); }
213
214     DPRINTF(IDL_FINISHED, ("sa2msghdr: 1\n"));
215     keyinfo->src = (struct sockaddr *)cp;
216     if (secassoc->src.sin6_len) {
217         bcopy((char *)&(secassoc->src), cp, secassoc->src.sin6_len);
218         ADVANCE(cp, secassoc->src.sin6_len);
219     } else {
220         bzero(cp, sizeof(struct sockaddr_in6));
221         ADVANCE(cp, sizeof(struct sockaddr_in6));
222     }
223     DPRINTF(IDL_FINISHED, ("sa2msghdr: 2\n"));
224
225     keyinfo->dst = (struct sockaddr *)&(secassoc->dst);
226     if (secassoc->dst.sin6_len) {
227         bcopy((char *)&(secassoc->dst), cp, secassoc->dst.sin6_len);
228         ADVANCE(cp, secassoc->dst.sin6_len);
229     } else {
230         bzero(cp, sizeof(struct sockaddr_in6));
231         ADVANCE(cp, sizeof(struct sockaddr_in6));
232     }
233     DPRINTF(IDL_FINISHED, ("sa2msghdr: 3\n"));
234
235     keyinfo->from = (struct sockaddr *)cp;
236     if (secassoc->from.sin6_len) {
237         bcopy((char *)&(secassoc->from), cp, secassoc->from.sin6_len);
238         ADVANCE(cp, secassoc->from.sin6_len);
239     } else {
240         bzero(cp, sizeof(struct sockaddr_in6));
241         ADVANCE(cp, sizeof(struct sockaddr_in6));
242     }
243     DPRINTF(IDL_FINISHED, ("sa2msghdr: 4\n"));
244
245     keyinfo->key = cp;
246     keyinfo->keylen = secassoc->keylen;
247     if (secassoc->keylen) {
248         bcopy((char *)(secassoc->key), cp, secassoc->keylen);
249         ADVANCE(cp, secassoc->keylen);
250     }
251
252     DPRINTF(IDL_FINISHED, ("sa2msghdr: 5\n"));
253     keyinfo->iv = cp;
254     keyinfo->ivlen = secassoc->ivlen;
255     if (secassoc->ivlen) {
256         bcopy((char *)(secassoc->iv), cp, secassoc->ivlen);
257         ADVANCE(cp, secassoc->ivlen);
258     }
259
260     DDO(IDL_FINISHED,printf("msgbuf(len=%d):\n",(char *)cp - (char *)km));
261     DDO(IDL_FINISHED,dump_buf((char *)km, (char *)cp - (char *)km));
262     DPRINTF(IDL_FINISHED, ("sa2msghdr: 6\n"));
```

4

```
263    return(0);
264  }
265
266
267  /*------------------------------------------------------------------
268   * key msghdr2secassoc():
269   *      Copy info from a key message buffer into an ipsec_assoc
270   *      structure
271   ------------------------------------------------------------------*
     /
272  int
273  key_msghdr2secassoc(secassoc, km, keyinfo)
274       struct ipsec assoc *secassoc;
275       struct key msghdr *km;
276       struct key_msgdata *keyinfo;
277  {
278    DPRINTF(IDL_GROSS_EVENT, ("Entering key_msghdr2secassoc\n"));
279
280    if ((km == 0) || (keyinfo == 0) || (secassoc == 0))
281      return(-1);
282
283    secassoc->len = sizeof(*secassoc);
284    secassoc->type = km->type;
285    secassoc->state = km->state;
286    secassoc->label = km->label;
287    secassoc->spi = km->spi;
288    secassoc->keylen = km->keylen;
289    secassoc->ivlen = km->ivlen;
290    secassoc->algorithm = km->algorithm;
291    secassoc->lifetype = km->lifetype;
292    secassoc->lifetime1 = km->lifetime1;
293    secassoc->lifetime2 = km->lifetime2;
294
295    if (keyinfo->src)
296      bcopy((char *)(keyinfo->src), (char *)&(secassoc->src),
297        keyinfo->src->sa_len);
298
299    if (keyinfo->dst)
300      bcopy((char *)(keyinfo->dst), (char *)&(secassoc->dst),
301        keyinfo->dst->sa_len);
302
303    if (keyinfo->from)
304      bcopy((char *)(keyinfo->from), (char *)&(secassoc->from),
305        keyinfo->from->sa_len);
306
307    /*
308     *  Make copies of key and iv
309     */
310    if (secassoc->ivlen) {
311      K Malloc(secassoc->iv, caddr_t, secassoc->ivlen);
312      if (secassoc->iv == 0) {
313        DPRINTF(IDL_CRITICAL,("msghdr2secassoc: can't allocate mem for
           iv\n"));
314        return(-1);
315      }
316      bcopy((char *)keyinfo->iv, (char *)secassoc->iv, secassoc->ivlen);
317    } else
318      secassoc->iv = NULL;
319
320    if (secassoc->keylen) {
321      K Malloc(secassoc->key, caddr_t, secassoc->keylen);
322      if (secassoc->key == 0) {
323        DPRINTF(IDL_CRITICAL,("msghdr2secassoc: can't allocate mem for
           key\n"));
324        if (secassoc->iv)
325      KFree(secassoc->iv);
326        return(-1);
```

5

```
327         }
328         bcopy((char *)keyinfo->key, (char *)secassoc->key,
            secassoc->keylen);
329       } else
330         secassoc->key = NULL;
331     return(0);
332  }
333
334
335  /*------------------------------------------------------------------
336   * addrpart_equal():
337   *       Determine if the address portion of two sockaddrs are equal.
338   *       Currently handles only AF_INET and AF_INET6 address families.
339    -----------------------------------------------------------------*
       /
340  int
341  addrpart_equal(sa1, sa2)
342       struct sockaddr *sa1;
343       struct sockaddr *sa2;
344  {
345
346     if ((sa1->sa_family == sa2->sa_family))
347       switch(sa1->sa_family) {
348       case AF_INET:
349         if (((struct sockaddr_in *)sa1)->sin_addr.s_addr ==
350         ((struct sockaddr_in *)sa2)->sin_addr.s_addr)
351       return(1);
352         break;
353       case AF_INET6:
354         if (IN6_ADDR_EQUAL(((struct sockaddr_in6 *)sa1)->sin6_addr,
355                  ((struct sockaddr_in6 *)sa2)->sin6_addr))
356       return(1);
357         break;
358       }
359     return(0);
360  }
361
362
363  /*------------------------------------------------------------------
364   * my_addr():
365   *       Determine if an address belongs to one of my configured
       interfaces.
366   *       Currently handles only AF_INET and AF_INET6 addresses.
367    -----------------------------------------------------------------*
       /
368  int
369  my_addr(sa)
370       struct sockaddr *sa;
371  {
372     extern struct in6_ifaddr *in6_ifaddr;
373     extern struct in_ifaddr *in_ifaddr;
374     struct in6_ifaddr *i6a = 0;
375     struct in_ifaddr *ia = 0;
376
377     switch(sa->sa_family) {
378     case AF_INET6:
379       for (i6a = in6_ifaddr; i6a; i6a = i6a->i6a_next) {
380         if (IN6_ADDR_EQUAL(((struct sockaddr_in6 *)sa)->sin6_addr,
381                  i6a->i6a_addr.sin6_addr))
382       return(1);
383       }
384       break;
385     case AF_INET:
386       for (ia = in_ifaddr; ia; ia = ia->ia_next) {
387         if (((struct sockaddr_in *)sa)->sin_addr.s_addr ==
388         ia->ia_addr.sin_addr.s_addr)
389       return(1);
```

6

```
390        }
391      break;
392      }
393    return(0);
394  }
395
396
397  /*-------------------------------------------------------------------
398   * key inittables():
399   *        Allocate space and initialize key engine tables
400   -------------------------------------------------------------------*
        /
401  void
402  key_inittables()
403  {
404    struct key_tblnode *keynode;
405    int i;
406
407    K_Malloc(keyregtable, struct key_registry *, sizeof(struct
       key_registry));
408    if (keyregtable == 0)
409      panic("key inittables");
410    bzero((char *)keyregtable, sizeof(struct key_registry));
411    K_Malloc(key_acquirelist, struct key_acquirelist *,
412        sizeof(struct key_acquirelist));
413    if (key_acquirelist == 0)
414      panic("key inittables");
415    bzero((char *)key_acquirelist, sizeof(struct key_acquirelist));
416    for (i = 0; i < KEYTBLSIZE; i++)
417      bzero((char *)&keytable[i], sizeof(struct key_tblnode));
418    for (i = 0; i < KEYALLOCTBLSIZE; i++)
419      bzero((char *)&keyalloctbl[i], sizeof(struct key_allocnode));
420    for (i = 0; i < SO2SPITBLSIZE; i++)
421      bzero((char *)&so2spitbl[i], sizeof(struct key_so2spinode));
422  }
423
424
425  /*-------------------------------------------------------------------
426   * key gethashval():
427   *        Determine keytable hash value.
428   -------------------------------------------------------------------*
        /
429  int
430  key_gethashval(buf, len, tblsize)
431        char *buf;
432        int len;
433        int tblsize;
434  {
435    int i, j = 0;
436
437    /*
438     * Todo: Use word size xor and check for alignment
439     *        and zero pad if necessary.  Need to also pick
440     *        a good hash function and table size.
441     */
442    if (len <= 0) {
443      DPRINTF(IDL_CRITICAL,("key_gethashval got bogus len!\n"));
444      return(-1);
445    }
446    for(i = 0; i < len; i++) {
447      j ^= (u_int8)(*(buf + i));
448    }
449    return (j % tblsize);
450  }
451
452
453  /*-------------------------------------------------------------------
```

7

```
454     * key createkey():
455     *         Create hash key for hash function
456     *         key is: type+src+dst if keytype = 1
457     *                 type+src+dst+spi if keytype = 0
458     *         Uses only the address portion of the src and dst sockaddrs to
459     *         form key.  Currently handles only AF INET and AF INET6 sockaddrs
460     ---------------------------------------------------------------------*
        /
461    int
462    key_createkey(buf, type, src, dst, spi, keytype)
463         char *buf;
464         u int type;
465         struct sockaddr *src;
466         struct sockaddr *dst;
467         u int32 spi;
468         u_int keytype;
469    {
470      char *cp, *p;
471
472      DPRINTF(IDL_FINISHED,("Entering key_createkey\n"));
473
474      if (!buf || !src || !dst)
475        return(-1);
476
477      cp = buf;
478      bcopy((char *)&type, cp, sizeof(type));
479      cp += sizeof(type);
480
481      /*
482       * Assume only IPv4 and IPv6 addresses.
483       */
484    #define ADDRPART(a) \
485        ((a)->sa family == AF INET6) ? \
486        (char *)&(((struct sockaddr in6 *)(a))->sin6 addr) : \
487        (char *)&(((struct sockaddr_in *)(a))->sin_addr)
488
489    #define ADDRSIZE(a) \
490        ((a)->sa family == AF INET6) ? sizeof(struct in_addr6) : \
491        sizeof(struct in_addr)
492
493      DPRINTF(IDL GROSS EVENT,("src addr:\n"));
494      DDO(IDL GROSS EVENT,dump smart sockaddr(src));
495      DPRINTF(IDL GROSS EVENT,("dst addr:\n"));
496      DDO(IDL_GROSS_EVENT,dump_smart_sockaddr(dst));
497
498      p = ADDRPART(src);
499      bcopy(p, cp, ADDRSIZE(src));
500      cp += ADDRSIZE(src);
501
502      p = ADDRPART(dst);
503      bcopy(p, cp, ADDRSIZE(dst));
504      cp += ADDRSIZE(dst);
505
506    #undef ADDRPART
507    #undef ADDRSIZE
508
509      if (keytype == 0) {
510        bcopy((char *)&spi, cp, sizeof(spi));
511        cp += sizeof(spi);
512      }
513
514      DPRINTF(IDL FINISHED,("hash key:\n"));
515      DDO(IDL FINISHED, dump_buf(buf, cp - buf));
516      return(cp - buf);
517    }
518
519
```

8

```
520  /*-------------------------------------------------------------------
521   * key sosearch():
522   *       Search the so2spi table for the security association allocated
      to
523   *       the socket.  Returns pointer to a struct key_so2spinode which
      can
524   *       be used to locate the security association entry in the
      keytable.
525   -------------------------------------------------------------------*
      /
526  struct key so2spinode *
527  key_sosearch(type, src, dst, so)
528       u int type;
529       struct sockaddr *src;
530       struct sockaddr *dst;
531       struct socket *so;
532  {
533    struct key_so2spinode *np = 0;
534
535    if (!(src && dst)) {
536      DPRINTF(IDL CRITICAL,("key_sosearch: got null src or dst
         pointer!\n"));
537      return(NULL);
538    }
539
540    for (np = so2spitbl[((u_int32)so) % SO2SPITBLSIZE].next; np; np = np->
       next) {
541      if ((so == np->socket) && (type == np->keynode->secassoc->type)
542      && addrpart equal(src,
543               (struct sockaddr *)&(np->keynode->secassoc->src))
544      && addrpart equal(dst,
545               (struct sockaddr *)&(np->keynode->secassoc->dst)))
546        return(np);
547    }
548    return(NULL);
549  }
550
551
552  /*-------------------------------------------------------------------
553   * key sodelete():
554   *       Delete entries from the so2spi table.
555   *          flag = 1   purge all entries
556   *          flag = 0   delete entries with socket pointer matching socket
557   -------------------------------------------------------------------*
      /
558  void
559  key_sodelete(socket, flag)
560       struct socket *socket;
561       int flag;
562  {
563    struct key so2spinode *prevnp, *np;
564    int s = splnet();
565
566    DPRINTF(IDL MAJOR_EVENT,("Entering keysodelete w/so-0x%x flag-%d\n",
       socket,flag));
567    if (flag) {
568      int i;
569
570      for (i = 0; i < SO2SPITBLSIZE; i++)
571        for(np = so2spitbl[i].next; np; np = np->next) {
572      KFree(np);
573        }
574      splx(s);
575      return;
576    }
577
578    prevnp = &so2spitbl[((u_int32)socket) % SO2SPITBLSIZE];
```

9

```
579        for(np = prevnp->next; np; np = np->next) {
580          if (np->socket == socket) {
581            struct socketlist *socklp, *prevsocklp;
582
583            (np->keynode->alloc_count)--;
584
585            /*
586             * If this socket maps to a unique secassoc,
587             * we go ahead and delete the secassoc, since it
588             * can no longer be allocated or used by any other
589             * socket.
590             */
591            if (np->keynode->secassoc->state & K_UNIQUE) {
592              if (key_delete(np->keynode->secassoc) != 0)
593                panic("key_sodelete");
594              np = prevnp;
595              continue;
596            }
597
598            /*
599             * We traverse the socketlist and remove the entry
600             * for this socket
601             */
602            DPRINTF(IDL_FINISHED,("keysodelete: deleting from socklist..."));
603            prevsocklp = np->keynode->solist;
604            for (socklp = prevsocklp->next; socklp; socklp = socklp->next) {
605              if (socklp->socket == socket) {
606                prevsocklp->next = socklp->next;
607                KFree(socklp);
608                break;
609              }
610              prevsocklp = socklp;
611            }
612            DPRINTF(IDL_FINISHED,("done\n"));
613            prevnp->next = np->next;
614            KFree(np);
615            np = prevnp;
616          }
617          prevnp = np;
618        }
619      splx(s);
620    }
621
622
623    /*------------------------------------------------------------------
624     * key deleteacquire():
625     *      Delete an entry from the key acquirelist
626     ------------------------------------------------------------------*
        /
627    void
628    key_deleteacquire(type, target)
629         u int type;
630         struct sockaddr *target;
631    {
632      struct key_acquirelist *ap, *prev;
633
634      prev = key acquirelist;
635      for(ap = key acquirelist->next; ap; ap = ap->next) {
636        if (addrpart equal(target, (struct sockaddr *)&(ap->target)) &&
637          (type == ap->type)) {
638          DPRINTF(IDL MAJOR EVENT,("Deleting entry from acquire list!\n"));
639          prev->next = ap->next;
640          KFree(ap);
641          ap = prev;
642        }
643        prev = ap;
644      }
```

10

```
645  }
646
647
648  /*-----------------------------------------------------------------------
649   * key search():
650   *        Search the key table for an entry with same type, src addr, dest
651   *        addr, and spi.  Returns a pointer to struct key_tblnode if found
652   *        else returns null.
653   *  ----------------------------------------------------------------------*
     /
654  struct key tblnode *
655  key_search(type, src, dst, spi, indx, prevkeynode)
656         u int type;
657         struct sockaddr *src;
658         struct sockaddr *dst;
659         u int32 spi;
660         int indx;
661         struct key_tblnode **prevkeynode;
662  {
663    struct key_tblnode *keynode, *prevnode;
664
665    if (indx > KEYTBLSIZE || indx < 0)
666      return (NULL);
667    if (!(&keytable[indx]))
668      return (NULL);
669
670  #define sec type keynode->secassoc->type
671  #define sec spi keynode->secassoc->spi
672  #define sec src keynode->secassoc->src
673  #define sec_dst keynode->secassoc->dst
674
675    prevnode = &keytable[indx];
676    for (keynode = keytable[indx].next; keynode; keynode = keynode->next)
       {
677      if ((type == sec type) && (spi == sec spi) &&
678      addrpart equal(src, (struct sockaddr *)&(sec src))
679      && addrpart_equal(dst, (struct sockaddr *)&(sec_dst)))
680        break;
681      prevnode = keynode;
682    }
683    *prevkeynode = prevnode;
684    return(keynode);
685  }
686
687
688  /*-----------------------------------------------------------------------
689   * key addnode():
690   *        Insert a key_tblnode entry into the key table.  Returns a
       pointer
691   *        to the newly created key tblnode.
692   *  ----------------------------------------------------------------------*
     /
693  struct key tblnode *
694  key_addnode(indx, secassoc)
695         int indx;
696         struct ipsec_assoc *secassoc;
697  {
698    struct key_tblnode *keynode;
699
700    DPRINTF(IDL GROSS EVENT,("Entering key addnode w/indx=%d
       secassoc=0x%x\n",indx, (u_int32)secassoc));
701
702    if (!(&keytable[indx]))
703      return(NULL);
704    if (!secassoc) {
705      panic("key_addnode: Someone passed in a null secassoc!\n");
706    }
```

11

```
707
708     K Malloc(keynode, struct key_tblnode *, sizeof(struct key_tblnode));
709     if (keynode == 0)
710        return(NULL);
711     bzero((char *)keynode, sizeof(struct key_tblnode));
712
713     K Malloc(keynode->solist, struct socketlist *, sizeof(struct
        socketlist));
714     if (keynode->solist == 0) {
715        KFree(keynode);
716        return(NULL);
717     }
718     bzero((char *)(keynode->solist), sizeof(struct socketlist));
719
720     keynode->secassoc = secassoc;
721     keynode->solist->next = NULL;
722     keynode->next = keytable[indx].next;
723     keytable[indx].next = keynode;
724     return(keynode);
725  }
726
727
728  /*-------------------------------------------------------------------------
729   * key add():
730   *       Add a new security association to the key table.  Caller is
731   *       responsible for allocating memory for the struct ipsec_assoc as
732   *       well as the buffer space for the key and iv.  Assumes the
        security
733   *       association passed in is well-formed.
734     -------------------------------------------------------------------------*
      /
735  int
736  key_add(secassoc)
737       struct ipsec_assoc *secassoc;
738  {
739     char buf[MAXHASHKEYLEN];
740     int len, indx;
741     int inbound = 0;
742     int outbound = 0;
743     struct key tblnode *keynode, *prevkeynode;
744     struct key_allocnode *np;
745     int s;
746
747     DPRINTF(IDL_GROSS_EVENT, ("Entering key_add w/secassoc=0x%x\n",
        secassoc));
748
749     if (!secassoc) {
750        panic("key_add: who the hell is passing me a null pointer");
751     }
752
753     /*
754      *  For storage purposes, the two esp modes are
755      *  treated the same.
756      */
757     if (secassoc->type == SS ENCRYPTION NETWORK)
758        secassoc->type = SS_ENCRYPTION_TRANSPORT;
759
760     /*
761      * Should we allow a null key to be inserted into the table ?
762      * or can we use null key to indicate some policy action...
763      */
764
765     /*
766      *  For esp using des-cbc or tripple-des we call
767      * des_set_odd_parity.
768      */
```

12

```
769     if (secassoc->key && (secassoc->type == SS ENCRYPTION TRANSPORT) &&
770         ((secassoc->algorithm == IPSEC ALGTYPE ESP DES CBC) ||
771          (secassoc->algorithm == IPSEC ALGTYPE_ESP_3DES)))
772       des_set_odd_parity(secassoc->key);
773
774     /*
775      *  Check if secassoc with same spi exists before adding
776      */
777     bzero((char *)&buf, sizeof(buf));
778     len = key_createkey((char *)&buf, secassoc->type,
779                     (struct sockaddr *)&(secassoc->src),
780                     (struct sockaddr *)&(secassoc->dst),
781                     secassoc->spi, 0);
782     indx = key gethashval((char *)&buf, len, KEYTBLSIZE);
783     DPRINTF(IDL GROSS EVENT,("keyadd: keytbl hash position=%d\n", indx));
784     keynode = key_search(secassoc->type, (struct sockaddr *)&(secassoc->
        src),
785                     (struct sockaddr *)&(secassoc->dst),
786                     secassoc->spi, indx, &prevkeynode);
787     if (keynode) {
788       DPRINTF(IDL_MAJOR_EVENT,("keyadd: secassoc already exists!\n"));
789       return(-2);
790     }
791
792     inbound = my addr((struct sockaddr *)&(secassoc->dst));
793     outbound = my addr((struct sockaddr *)&(secassoc->src));
794     DPRINTF(IDL_FINISHED,("inbound=%d outbound=%d\n", inbound, outbound));
795
796     /*
797      * We allocate mem for an allocation entry if needed.
798      * This is done here instead of in the allocaton code
799      * segment so that we can easily recover/cleanup from a
800      * memory allocation error.
801      */
802     if (outbound || (!inbound && !outbound)) {
803       K Malloc(np, struct key_allocnode *, sizeof(struct key_allocnode));
804       if (np == 0) {
805         DPRINTF(IDL_CRITICAL,("keyadd: can't allocate allocnode!\n"));
806         return(-1);
807       }
808     }
809
810     s = splnet();
811
812     if ((keynode = key addnode(indx, secassoc)) == NULL) {
813       DPRINTF(IDL_CRITICAL,("keyadd: key_addnode failed!\n"));
814       if (np)
815         KFree(np);
816       splx(s);
817       return(-1);
818     }
819     DPRINTF(IDL EVENT,("Added new keynode:\n"));
820     DDO(IDL GROSS EVENT, dump keytblnode(keynode));
821     DDO(IDL_GROSS_EVENT, dump_secassoc(keynode->secassoc));
822
823     /*
824      *  We add an entry to the allocation table for
825      *  this secassoc if the interfaces are up and
826      *  the secassoc is outbound.  In the case
827      *  where the interfaces are not up, we go ahead
828      *  and do it anyways.  This wastes an allocation
829      *  entry if the secassoc later turned out to be
830      *  inbound when the interfaces are ifconfig up.
831      */
832     if (outbound || (!inbound && !outbound)) {
833       len = key createkey((char *)&buf, secassoc->type,
834                     (struct sockaddr *)&(secassoc->src),
```

13

```
835                   (struct sockaddr *)&(secassoc->dst),
836                 0, 1);
837        indx = key gethashval((char *)&buf, len, KEYALLOCTBLSIZE);
838        DPRINTF(IDL_GROSS_EVENT,("keyadd: keyalloc hash position=%d\n",
           indx));
839        np->keynode = keynode;
840        np->next = keyalloctbl[indx].next;
841        keyalloctbl[indx].next = np;
842      }
843      if (inbound)
844        secassoc->state |= K_INBOUND;
845      if (outbound)
846        secassoc->state |= K_OUTBOUND;
847
848      key deleteacquire(secassoc->type, (struct sockaddr
         *)&(secassoc->dst));
849
850      splx(s);
851      return 0;
852    }
853
854
855    /*----------------------------------------------------------------
856     * key get():
857     *      Get a security association from the key table.
858     -----------------------------------------------------------------*
           /
859    int
860    key_get(type, src, dst, spi, secassoc)
861         u int type;
862         struct sockaddr *src;
863         struct sockaddr *dst;
864         u int32 spi;
865         struct ipsec_assoc **secassoc;
866    {
867      char buf[MAXHASHKEYLEN];
868      struct key tblnode *keynode, *prevkeynode;
869      int len, indx;
870
871      /*
872       *  For storage purposes, the two esp modes are
873       *  treated the same.
874       */
875      if (type == SS ENCRYPTION NETWORK)
876        type = SS_ENCRYPTION_TRANSPORT;
877
878      bzero(&buf, sizeof(buf));
879      *secassoc = NULL;
880      len = key createkey((char *)&buf, type, src, dst, spi, 0);
881      indx = key gethashval((char *)&buf, len, KEYTBLSIZE);
882      DPRINTF(IDL GROSS EVENT,("keyget: indx=%d\n",indx));
883      keynode = key search(type, src, dst, spi, indx, &prevkeynode);
884      if (keynode) {
885        DPRINTF(IDL_EVENT,("keyget: found it! keynode=0x%x",keynode));
886        *secassoc = keynode->secassoc;
887        return(0);
888      } else
889        return(-1);   /* Not found */
890    }
891
892
893    /*----------------------------------------------------------------
894     * key dump():
895     *      Dump all valid entries in the keytable to a pf key socket.  Each
896     *      security associaiton is sent one at a time in a pf_key message.
           A
897     *      message with seqno = 0 signifies the end of the dump
```

14

```
897    transaction.
898    -------------------------------------------------------------------*
       /
899    int
900    key_dump(so)
901         struct socket *so;
902    {
903      int len, i;
904      int seq = 1;
905      struct mbuf *m;
906      struct key msgdata keyinfo;
907      struct key msghdr *km;
908      struct key tblnode *keynode;
909      extern struct sockaddr key src;
910      extern struct sockaddr key_dst;
911
912      /*
913       * Routine to dump the key table to a routing socket
914       * Use for debugging only!
915       */
916
917      DPRINTF(IDL_GROSS_EVENT,("Entering key_dump()"));
918      /*
919       * We need to speed this up later.  Fortunately, key dump
920       * messages are not sent often.
921       */
922      for (i = 0; i < KEYTBLSIZE; i++) {
923        for (keynode = keytable[i].next; keynode; keynode = keynode->next) {
924
925          /*
926           * We exclude dead/larval/zombie security associations for now
927           * but it may be useful to also send these up for debugging
             purposes
928           */
929          if (keynode->secassoc->state & (K_DEAD | K_LARVAL | K_ZOMBIE))
930        continue;
931
932          len = (sizeof(struct key msghdr) +
933              ROUNDUP(keynode->secassoc->src.sin6 len) +
934              ROUNDUP(keynode->secassoc->dst.sin6 len) +
935              ROUNDUP(keynode->secassoc->from.sin6_len) +
936              ROUNDUP(keynode->secassoc->keylen) +
937              ROUNDUP(keynode->secassoc->ivlen));
938          K Malloc(km, struct key_msghdr *, len);
939          if (km == 0)
940        return(ENOBUFS);
941          if (key secassoc2msghdr(keynode->secassoc, km, &keyinfo) != 0)
942        panic("key dump");
943          km->key msglen = len;
944          km->key msgvers = KEY VERSION;
945          km->key msgtype = KEY DUMP;
946          km->key pid = curproc->p_pid;
947          km->key seq = seq++;
948          km->key errno = 0;
949          MGETHDR(m, M WAIT, MT DATA);
950          m->m len = m->m_pkthdr.len = 0;
951          m->m next = 0;
952          m->m nextpkt = 0;
953          m->m pkthdr.rcvif = 0;
954          m copyback(m, 0, len, (caddr_t)km);
955          KFree(km);
956          if (sbappendaddr(&so->so_rcv, &key_src, m, (struct mbuf *)0) == 0)
957        m_free(m);
958          else
959        sorwakeup(so);
960        }
961    }
```

15

```
962    K Malloc(km, struct key_msghdr *, sizeof(struct key_msghdr));
963    if (km == 0)
964      return(ENOBUFS);
965    bzero((char *)km, sizeof(struct key msghdr));
966    km->key msglen = sizeof(struct key_msghdr);
967    km->key msgvers = KEY VERSION;
968    km->key msgtype = KEY DUMP;
969    km->key pid = curproc->p_pid;
970    km->key seq = 0;
971    km->key errno = 0;
972    MGETHDR(m, M WAIT, MT DATA);
973    m->m len = m->m_pkthdr.len = 0;
974    m->m next = 0;
975    m->m nextpkt = 0;
976    m->m pkthdr.rcvif = 0;
977    m copyback(m, 0, km->key_msglen, (caddr_t)km);
978    KFree(km);
979    if (sbappendaddr(&so->so_rcv, &key_src, m, (struct mbuf *)0) == 0)
980      m free(m);
981    else
982      sorwakeup(so);
983    DPRINTF(IDL_GROSS_EVENT,("Leaving key_dump()\n"));
984    return(0);
985  }
986
987  /*-------------------------------------------------------------------
988   * key delete():
989   *      Delete a security association from the key table.
990   -------------------------------------------------------------------*
     /
991  int
992  key_delete(secassoc)
993        struct ipsec_assoc *secassoc;
994  {
995    char buf[MAXHASHKEYLEN];
996    int len, indx;
997    struct key tblnode *keynode = 0;
998    struct key tblnode *prevkeynode = 0;
999    struct socketlist *socklp, *deadsocklp;
1000   struct key so2spinode *np, *prevnp;
1001   struct key_allocnode *ap, *prevap;
1002   int s;
1003
1004   DPRINTF(IDL_GROSS_EVENT,("Entering key_delete w/secassoc=0x%x\n",
       secassoc));
1005
1006   if (secassoc->type == SS ENCRYPTION NETWORK)
1007     secassoc->type = SS_ENCRYPTION_TRANSPORT;
1008
1009   bzero((char *)&buf, sizeof(buf));
1010   len = key_createkey((char *)&buf, secassoc->type,
1011              (struct sockaddr *)&(secassoc->src),
1012              (struct sockaddr *)&(secassoc->dst),
1013              secassoc->spi, 0);
1014   indx = key gethashval((char *)&buf, len, KEYTBLSIZE);
1015   DPRINTF(IDL_GROSS_EVENT,("keydelete: keytbl hash position=%d\n",
       indx));
1016   keynode = key_search(secassoc->type, (struct sockaddr *)&(secassoc->
       src),
1017              (struct sockaddr *)&(secassoc->dst),
1018              secassoc->spi, indx, &prevkeynode);
1019
1020   if (keynode) {
1021     s = splnet();
1022     DPRINTF(IDL EVENT,("keydelete: found keynode to delete\n"));
1023     keynode->secassoc->state |= K_DEAD;
1024
```

16

```
1025        if (keynode->ref count > 0) {
1026            DPRINTF(IDL MAJOR EVENT,("keydelete: secassoc still held, marking
                for deletion only!\n"));
1027            splx(s);
1028            return(0);
1029        }
1030
1031        prevkeynode->next = keynode->next;
1032
1033        /*
1034         *  Walk the socketlist and delete the
1035         *  entries mapping sockets to this secassoc
1036         *  from the so2spi table.
1037         */
1038        DPRINTF(IDL GROSS EVENT,("keydelete: deleting socklist..."));
1039        for(socklp = keynode->solist->next; socklp; ) {
1040            prevnp = &so2spitbl[((u int32)(socklp->socket)) % SO2SPITBLSIZE];
1041            for(np = prevnp->next; np; np = np->next) {
1042        if ((np->socket == socklp->socket) && (np->keynode == keynode)) {
1043            prevnp->next = np->next;
1044            KFree(np);
1045            break;
1046        }
1047        prevnp = np;
1048            }
1049            deadsocklp = socklp;
1050            socklp = socklp->next;
1051            KFree(deadsocklp);
1052        }
1053        DPRINTF(IDL_GROSS_EVENT,("done\n"));
1054        /*
1055         * If an allocation entry exist for this
1056         * secassoc, delete it.
1057         */
1058        bzero((char *)&buf, sizeof(buf));
1059        len = key createkey((char *)&buf, secassoc->type,
1060                    (struct sockaddr *)&(secassoc->src),
1061                    (struct sockaddr *)&(secassoc->dst),
1062                    0, 1);
1063        indx = key gethashval((char *)&buf, len, KEYALLOCTBLSIZE);
1064        DPRINTF(IDL_GROSS_EVENT,("keydelete: alloctbl hash position=%d\n",
            indx));
1065        prevap = &keyalloctbl[indx];
1066        for (ap = prevap->next; ap; ap = ap->next) {
1067            if (ap->keynode == keynode) {
1068        prevap->next = ap->next;
1069        KFree(ap);
1070        break;
1071            }
1072            prevap = ap;
1073        }
1074
1075        if (keynode->secassoc->iv)
1076            KFree(keynode->secassoc->iv);
1077        if (keynode->secassoc->key)
1078            KFree(keynode->secassoc->key);
1079        KFree(keynode->secassoc);
1080        if (keynode->solist)
1081            KFree(keynode->solist);
1082        KFree(keynode);
1083        splx(s);
1084        return(0);
1085    }
1086    return(-1);
1087 }
1088
1089
```

```
1090   /*------------------------------------------------------------------
1091    * key flush():
1092    *       Delete all entries from the key table.
1093        ------------------------------------------------------------------*
       /
1094   void
1095   key_flush(void)
1096   {
1097     struct key_tblnode *keynode;
1098     int i;
1099
1100     /*
1101      * This is slow, but simple.
1102      */
1103     DPRINTF(IDL_FINISHED,("Flushing key table..."));
1104     for (i = 0; i < KEYTBLSIZE; i++) {
1105       while (keynode = keytable[i].next)
1106         if (key delete(keynode->secassoc) != 0)
1107       panic("key_flush");
1108     }
1109     DPRINTF(IDL_FINISHED,("done\n"));
1110   }
1111
1112
1113   /*------------------------------------------------------------------
1114    * key getspi():
1115    *       Get a unique spi value for a key management daemon/program.  The
1116    *       spi value, once assigned, cannot be assigned again.
1117        ------------------------------------------------------------------*
       /
1118   int
1119   key_getspi(type, src, dst, spi)
1120        u int type;
1121        struct sockaddr *src;
1122        struct sockaddr *dst;
1123        u_int32 *spi;
1124   {
1125     struct ipsec assoc *secassoc;
1126     struct key_tblnode *keynode, *prevkeynode;
1127     int count, done, len, indx;
1128     int maxcount = 1000;
1129     u int32 val;
1130     char buf[MAXHASHKEYLEN];
1131     int s;
1132
1133
1134     DPRINTF(IDL MAJOR EVENT,("Entering getspi w/type=%d\n",type));
1135     if (!(src && dst))
1136       return(-1);
1137
1138     /*
1139      *  For storage purposes, the two esp modes are
1140      *  treated the same.
1141      */
1142     if (type == SS ENCRYPTION NETWORK)
1143       type = SS_ENCRYPTION_TRANSPORT;
1144
1145     done = count = 0;
1146     do {
1147       count++;
1148       /*
1149        *  Currently, valid spi values are 32 bits wide except for
1150        *  the value of zero.  This need to change to take into
1151        *  account more restrictive spi ranges.
1152        *
1153        *  TODO:  Kebe says to allow key mgnt daemon to specify range
```

18

```
1154        *           of valid spi to get.
1155        */
1156       val = random();
1157       DPRINTF(IDL_FINISHED,("%u ",val));
1158       if (val) {
1159         DPRINTF(IDL_FINISHED,("\n"));
1160         bzero(&buf, sizeof(buf));
1161         len = key_createkey((char *)&buf, type, src, dst, val, 0);
1162         indx = key_gethashval((char *)&buf, len, KEYTBLSIZE);
1163         if(!key_search(type, src, dst, val, indx, &prevkeynode)) {
1164       s = splnet();
1165       K_Malloc(secassoc, struct ipsec_assoc *, sizeof(struct
           ipsec_assoc));
1166       if (secassoc == 0) {
1167         DPRINTF(IDL_CRITICAL,("key_getspi: can't allocate memory\n"));
1168         splx(s);
1169         return(-1);
1170       }
1171       bzero((char *)secassoc, sizeof(struct ipsec_assoc));
1172
1173       DPRINTF(IDL_FINISHED,("getspi: indx=%d\n",indx));
1174       secassoc->len = sizeof(struct ipsec_assoc);
1175       secassoc->type = type;
1176       secassoc->spi = val;
1177       secassoc->state |= K_LARVAL;
1178       if (my_addr((struct sockaddr *)&(secassoc->dst)))
1179         secassoc->state |= K_INBOUND;
1180       if (my_addr((struct sockaddr *)&(secassoc->src)))
1181         secassoc->state |= K_OUTBOUND;
1182
1183       bcopy((char *)src, (char *)&(secassoc->src), src->sa_len);
1184       bcopy((char *)dst, (char *)&(secassoc->dst), dst->sa_len);
1185       secassoc->from.sin6_family = AF_INET6;
1186       secassoc->from.sin6_len = sizeof(struct sockaddr_in6);
1187
1188       /*
1189        * We need to add code to age these larval key table
1190        * entries so they don't linger forever waiting for
1191        * a KEY_UPDATE message that may not come for various
1192        * reasons.  This is another task that key_reaper can
1193        * do once we have it coded.
1194        */
1195       secassoc->lifetime1 = time.tv_sec + maxlarvallifetime;
1196
1197       if (!(keynode = key_addnode(indx, secassoc))) {
1198         DPRINTF(IDL_CRITICAL,("key_getspi: can't add node\n"));
1199         splx(s);
1200         return(-1);
1201       }
1202       DPRINTF(IDL_FINISHED,("key_getspi: added node 0x%x\n",keynode));
1203       done++;
1204       splx(s);
1205         }
1206       }
1207     } while ((count < maxcount) && !done);
1208     DPRINTF(IDL_FINISHED,("getspi returns
          w/spi=%u,count=%d\n",val,count));
1209     if (done) {
1210       *spi = val;
1211       return(0);
1212     } else {
1213       *spi = 0;
1214       return(-1);
1215     }
1216 }
1217
1218
```

19

DEF00007960

```
1219   /*-----------------------------------------------------------------
1220    * key update():
1221    *       Update a keytable entry that has an spi value assigned but is
1222    *       incomplete (e.g. no key/iv).
1223    -----------------------------------------------------------------*
       /
1224   int
1225   key_update(secassoc)
1226        struct ipsec_assoc *secassoc;
1227   {
1228     struct key tblnode *keynode, *prevkeynode;
1229     struct key allocnode *np = 0;
1230     u int8 newstate;
1231     int len, indx, inbound, outbound;
1232     char buf[MAXHASHKEYLEN];
1233     int s;
1234
1235     /*
1236      *  For storage purposes, the two esp modes are
1237      *  treated the same.
1238      */
1239     if (secassoc->type == SS ENCRYPTION NETWORK)
1240       secassoc->type = SS_ENCRYPTION_TRANSPORT;
1241
1242     bzero(&buf, sizeof(buf));
1243     len = key_createkey((char *)&buf, secassoc->type,
1244                  (struct sockaddr *)&(secassoc->src),
1245                  (struct sockaddr *)&(secassoc->dst),
1246                  secassoc->spi, 0);
1247     indx = key gethashval((char *)&buf, len, KEYTBLSIZE);
1248     if(!(keynode = key search(secassoc->type,
1249                  (struct sockaddr *)&(secassoc->src),
1250                  (struct sockaddr *)&(secassoc->dst),
1251                  secassoc->spi, indx, &prevkeynode))) {
1252       return(ESRCH);
1253     }
1254     if (keynode->secassoc->state & K_DEAD)
1255       return(ESRCH);
1256
1257     /* Should we also restrict updating of only LARVAL entries ? */
1258
1259     s = splnet();
1260
1261     inbound = my addr((struct sockaddr *)&(secassoc->dst));
1262     outbound = my_addr((struct sockaddr *)&(secassoc->src));
1263
1264     newstate = keynode->secassoc->state;
1265     newstate &= ~K_LARVAL;
1266     if (inbound)
1267       newstate |= K_INBOUND;
1268     if (outbound)
1269       newstate |= K_OUTBOUND;
1270
1271     if (outbound || (!inbound && !outbound)) {
1272       K Malloc(np, struct key_allocnode *, sizeof(struct key_allocnode));
1273       if (np == 0) {
1274         DPRINTF(IDL_CRITICAL,("keyupdate: can't allocate allocnode!\n"));
1275         splx(s);
1276         return(ENOBUFS);
1277       }
1278     }
1279
1280     /*
1281      *  We now copy the secassoc over. We don't need to copy
1282      *  the key and iv into new buffers since the calling routine
1283      *  does that already.
1284      */
```

20

```
1285
1286      *(keynode->secassoc) = *secassoc;
1287      keynode->secassoc->state = newstate;
1288
1289      /*
1290       * Should we allow a null key to be inserted into the table ?
1291       * or can we use null key to indicate some policy action...
1292       */
1293
1294      if (keynode->secassoc->key &&
1295           (keynode->secassoc->type == SS_ENCRYPTION TRANSPORT) &&
1296           ((keynode->secassoc->algorithm == IPSEC ALGTYPE ESP DES_CBC) ||
1297        (keynode->secassoc->algorithm == IPSEC ALGTYPE_ESP_3DES)))
1298         des_set_odd_parity(keynode->secassoc->key);
1299
1300      /*
1301       *  We now add an entry to the allocation table for this
1302       *  updated key table entry.
1303       */
1304      if (outbound || (!inbound && !outbound)) {
1305        len = key createkey((char *)&buf, secassoc->type,
1306                (struct sockaddr *)&(secassoc->src),
1307                (struct sockaddr *)&(secassoc->dst),
1308                0, 1);
1309        indx = key gethashval((char *)&buf, len, KEYALLOCTBLSIZE);
1310        DPRINTF(IDL_FINISHED,("keyupdate: keyalloc hash position=%d\n",
                 indx));
1311        np->keynode = keynode;
1312        np->next = keyalloctbl[indx].next;
1313        keyalloctbl[indx].next = np;
1314      }
1315
1316      key deleteacquire(secassoc->type, (struct sockaddr
          *)&(secassoc->dst));
1317
1318      splx(s);
1319      return(0);
1320  }
1321
1322  /*-----------------------------------------------------------------
1323   * key register():
1324   *       Register a socket as one capable of acquiring security
          associations
1325   *       for the kernel.
1326      ---------------------------------------------------------------*
       /
1327  int
1328  key_register(socket, type)
1329       struct socket *socket;
1330       u_int type;
1331  {
1332    struct key registry *p, *new;
1333    int s = splnet();
1334
1335    DPRINTF(IDL MAJOR_EVENT,("Entering key_register w/so=0x%x,type=%d\n",
          socket,type));
1336
1337    if (!(keyregtable && socket))
1338      panic("key_register");
1339
1340    /*
1341     * Make sure entry is not already in table
1342     */
1343    for(p = keyregtable->next; p; p = p->next) {
1344      if ((p->type == type) && (p->socket == socket)) {
1345        splx(s);
1346        return(EEXIST);
```

21

```
1347        }
1348      }
1349
1350      K Malloc(new, struct key_registry *, sizeof(struct key_registry));
1351      if (new == 0) {
1352        splx(s);
1353        return(ENOBUFS);
1354      }
1355      new->type = type;
1356      new->socket = socket;
1357      new->next = keyregtable->next;
1358      keyregtable->next = new;
1359      splx(s);
1360      return(0);
1361   }
1362
1363   /*-----------------------------------------------------------------------
1364    * key unregister():
1365    *      Delete entries from the registry list.
1366    *          allflag = 1 : delete all entries with matching socket
1367    *          allflag = 0 : delete only the entry matching socket and type
1368      ----------------------------------------------------------------------*
       /
1369   void
1370   key_unregister(socket, type, allflag)
1371        struct socket *socket;
1372        u int type;
1373        int allflag;
1374   {
1375      struct key registry *p, *prev;
1376      int s = splnet();
1377
1378      DPRINTF(IDL MAJOR EVENT,("Entering key unregister
          w/so=0x%x,type=%d,flag=%d\n",socket, type, allflag));
1379
1380      if (!(keyregtable && socket))
1381        panic("key register");
1382      prev = keyregtable;
1383      for(p = keyregtable->next; p; p = p->next) {
1384        if ((allflag && (p->socket == socket)) ||
1385        ((p->type == type) && (p->socket == socket))) {
1386          prev->next = p->next;
1387          KFree(p);
1388          p = prev;
1389        }
1390        prev = p;
1391      }
1392      splx(s);
1393   }
1394
1395
1396   /*--------------------------------------------------------------------
1397    * key acquire():
1398    *      Send a key acquire message to all registered key mgnt daemons
1399    *      capable of acquire security association of type type.
1400    *
1401    *      Return: 0 if succesfully called key mgnt. daemon(s)
1402    *              -1 if not successfull.
1403      -------------------------------------------------------------------*
       /
1404   int
1405   key_acquire(type, src, dst)
1406        u int type;
1407        struct sockaddr *src;
1408        struct sockaddr *dst;
1409   {
1410      struct key_registry *p;
```

22

```
1411        struct key acquirelist *ap, *prevap;
1412        int success = 0, created = 0;
1413        struct socket *last = 0;
1414        struct mbuf *m = 0;
1415        u int etype;
1416        extern struct sockaddr key_src;
1417
1418        DPRINTF(IDL_MAJOR_EVENT,("Entering key_acquire()\n"));
1419
1420        if (!keyregtable || !src || !dst)
1421          return (-1);
1422
1423        /*
1424         * We first check the acquirelist to see if a key_acquire
1425         * message has been sent for this destination.
1426         */
1427        etype = type;
1428        if (etype == SS ENCRYPTION NETWORK)
1429          etype = SS ENCRYPTION TRANSPORT;
1430        prevap = key acquirelist;
1431        for(ap = key acquirelist->next; ap; ap = ap->next) {
1432          if (addrpart equal(dst, (struct sockaddr *)&(ap->target)) &&
1433          (etype == ap->type)) {
1434            DPRINTF(IDL MAJOR EVENT,("acquire message previously sent!\n"));
1435            if (ap->expiretime < time.tv sec) {
1436            DPRINTF(IDL_MAJOR_EVENT,("acquire message has expired!\n"));
1437            ap->count = 0;
1438            break;
1439            }
1440            if (ap->count < maxkeyacquire) {
1441            DPRINTF(IDL MAJOR_EVENT,("max acquire messages not yet
            exceeded!\n"));
1442            break;
1443            }
1444            return(0);
1445          } else if (ap->expiretime < time.tv_sec) {
1446            /*
1447             *  Since we're already looking at the list, we may as
1448             *  well delete expired entries as we scan through the list.
1449             *  This should really be done by a function like key reaper()
1450             *  but until we code key_reaper(), this is a quick and dirty
1451             *  hack.
1452             */
1453            DPRINTF(IDL_MAJOR_EVENT,("found an expired entry...deleting
            it!\n"));
1454            prevap->next = ap->next;
1455            KFree(ap);
1456            ap = prevap;
1457          }
1458          prevap = ap;
1459        }
1460
1461        /*
1462         * Scan registry and send KEY ACQUIRE message to
1463         * appropriate key management daemons.
1464         */
1465        for(p = keyregtable->next; p; p = p->next) {
1466          if (p->type != type)
1467            continue;
1468
1469          if (!created) {
1470            struct key_msghdr *km;
1471            int len;
1472
1473            len = sizeof(struct key_msghdr) + ROUNDUP(src->sa_len) +
1474          ROUNDUP(dst->sa len);
1475            K_Malloc(km, struct key_msghdr *, len);
```

23

```
1476        if (km == 0) {
1477      DPRINTF(IDL_CRITICAL,("key_acquire: no memory\n"));
1478      return(-1);
1479        }
1480        DPRINTF(IDL_FINISHED,("key_acquire/created: 1\n"));
1481        bzero((char *)km, len);
1482        km->key_msglen = len;
1483        km->key_msgvers = KEY_VERSION;
1484        km->key_msgtype = KEY_ACQUIRE;
1485        km->type = type;
1486        DPRINTF(IDL_FINISHED,("key_acquire/created: 2\n"));
1487        /*
1488         * This is inefficient and slow.
1489         */
1490
1491        /*
1492         * We zero out sin_zero here for AF_INET addresses because
1493         * ip_output() currently does not do it for performance reasons.
1494         */
1495        if (src->sa_family == AF_INET)
1496      bzero((char *)(&((struct sockaddr_in *)src)->sin_zero),
1497            sizeof(((struct sockaddr_in *)src)->sin_zero));
1498        if (dst->sa_family == AF_INET)
1499      bzero((char *)(&((struct sockaddr_in *)dst)->sin_zero),
1500            sizeof(((struct sockaddr_in *)dst)->sin_zero));
1501
1502        bcopy((char *)src, (char *)(km + 1), src->sa_len);
1503        bcopy((char *)dst, (char *)((int)(km + 1) + ROUNDUP(src->sa_len)),
1504          dst->sa_len);
1505        DPRINTF(IDL_FINISHED,("key_acquire/created: 3\n"));
1506        MGETHDR(m, M_WAIT, MT_DATA);
1507        m->m_len = m->m_pkthdr.len = 0;
1508        m->m_next = 0;
1509        m->m_nextpkt = 0;
1510        m->m_pkthdr.rcvif = 0;
1511        m_copyback(m, 0, len, (caddr_t)km);
1512        KFree(km);
1513        DPRINTF(IDL_FINISHED,("key_acquire/created: 4\n"));
1514        DDO(IDL_FINISHED,dump_mchain(m));
1515        created++;
1516        }
1517        if (last) {
1518      struct mbuf *n;
1519      if (n = m_copy(m, 0, (int)M_COPYALL)) {
1520        if (sbappendaddr(&last->so_rcv, &key_src, n, (struct mbuf *)0) == 0)
1521          m_freem(n);
1522        else {
1523          sorwakeup(last);
1524          success++;
1525        }
1526      }
1527      DPRINTF(IDL_FINISHED,("key_acquire/last: 1\n"));
1528        }
1529      last = p->socket;
1530    }
1531    if (last) {
1532      if (sbappendaddr(&last->so_rcv, &key_src, m, (struct mbuf *)0) == 0)
1533        m_freem(m);
1534      else {
1535        sorwakeup(last);
1536        success++;
1537      }
1538      DPRINTF(IDL_FINISHED,("key_acquire/last: 2\n"));
1539    } else
1540      m_freem(m);
1541
1542    /*
```

24

```
1543        *  Update the acquirelist
1544        */
1545       if (success) {
1546         if (!ap) {
1547           DPRINTF(IDL MAJOR EVENT,("Adding new entry in acquirelist\n"));
1548           K Malloc(ap, struct key_acquirelist *, sizeof(struct
              key acquirelist));
1549           if (ap == 0)
1550         return(success ? 0 : -1);
1551           bzero((char *)ap, sizeof(struct key acquirelist));
1552           bcopy((char *)dst, (char *)&(ap->target), dst->sa_len);
1553           ap->type = etype;
1554           ap->next = key acquirelist->next;
1555           key_acquirelist->next = ap;
1556         }
1557         DPRINTF(IDL_EVENT,("Updating acquire counter and expiration
              time\n"));
1558         ap->count++;
1559         ap->expiretime = time.tv_sec + maxacquiretime;
1560       }
1561       DPRINTF(IDL MAJOR EVENT,("key_acquire: done! success=%d\n",success));
1562       return(success ? 0 : -1);
1563   }
1564
1565   /*-----------------------------------------------------------------
1566    * key alloc():
1567    *       Allocate a security association to a socket.  A socket
           requesting
1568    *       unique keying (per-socket keying) is assigned a security
           assocation
1569    *       exclusively for its use.  Sockets not requiring unique keying
           are
1570    *       assigned the first security association which may or may not be
1571    *       used by another socket.
1572    -----------------------------------------------------------------*
           /
1573   int
1574   key_alloc(type, src, dst, socket, unique_key, keynodep)
1575         u int type;
1576         struct sockaddr *src;
1577         struct sockaddr *dst;
1578         struct socket *socket;
1579         u int  unique key;
1580         struct key_tblnode **keynodep;
1581   {
1582     struct key tblnode *keynode;
1583     char buf[MAXHASHKEYLEN];
1584     struct key allocnode *np, *prevnp;
1585     struct key_so2spinode *newnp;
1586     int len;
1587     int indx;
1588
1589     DPRINTF(IDL GROSS EVENT,("Entering key_alloc w/type=%u!\n",type));
1590     if (!(src && dst)) {
1591       DPRINTF(IDL_CRITICAL,("key_alloc: received null src or dst!\n"));
1592       return(-1);
1593     }
1594
1595     /*
1596      *  We treat esp-transport mode and esp-tunnel mode
1597      *  as a single type in the keytable.
1598      */
1599     if (type == SS ENCRYPTION NETWORK)
1600       type = SS_ENCRYPTION_TRANSPORT;
1601
1602     /*
1603      * Search key allocation table
```

25

```
1604      */
1605      bzero((char *)&buf, sizeof(buf));
1606      len = key_createkey((char *)&buf, type, src, dst, 0, 1);
1607      indx = key_gethashval((char *)&buf, len, KEYALLOCTBLSIZE);
1608
1609  #define np_type np->keynode->secassoc->type
1610  #define np_state np->keynode->secassoc->state
1611  #define np_src (struct sockaddr *)&(np->keynode->secassoc->src)
1612  #define np_dst (struct sockaddr *)&(np->keynode->secassoc->dst)
1613
1614      prevnp = &keyalloctbl[indx];
1615      for (np = keyalloctbl[indx].next; np; np = np->next) {
1616        if ((type == np_type) && addrpart_equal(src, np_src) &&
1617        addrpart_equal(dst, np_dst) &&
1618        !(np_state & (K_LARVAL | K_DEAD | K_UNIQUE))) {
1619          if (!(unique_key))
1620        break;
1621          if (!(np_state & K_USED))
1622        break;
1623        }
1624        prevnp = np;
1625      }
1626
1627      if (np) {
1628        struct key_so2spinode *newp;
1629        struct socketlist *newsp;
1630        int s = splnet();
1631
1632        DPRINTF(IDL_MAJOR_EVENT,("key_alloc: found node to allocate\n"));
1633        keynode = np->keynode;
1634
1635        K_Malloc(newp, struct key_so2spinode *, sizeof(struct
1636        key_so2spinode));
1636        if (newp == 0) {
1637          DPRINTF(IDL_CRITICAL,("key_alloc: Can't alloc mem for so2spi
1637          node!\n"));
1638          splx(s);
1639          return(ENOBUFS);
1640        }
1641        K_Malloc(newsp, struct socketlist *, sizeof(struct socketlist));
1642        if (newsp == 0) {
1643          DPRINTF(IDL_CRITICAL,("key_alloc: Can't alloc mem for
1643          socketlist!\n"));
1644          if (newp)
1645        KFree(newp);
1646          splx(s);
1647          return(ENOBUFS);
1648        }
1649
1650        /*
1651         * Add a hash entry into the so2spi table to
1652         * map socket to allocated secassoc.
1653         */
1654        DPRINTF(IDL_GROSS_EVENT,("key_alloc: adding entry to so2spi
1654        table..."));
1655        newp->keynode = keynode;
1656        newp->socket = socket;
1657        newp->next = so2spitbl[((u_int32)socket) % SO2SPITBLSIZE].next;
1658        so2spitbl[((u_int32)socket) % SO2SPITBLSIZE].next = newp;
1659        DPRINTF(IDL_GROSS_EVENT,("done\n"));
1660
1661        if (unique_key) {
1662          /*
1663           * Need to remove the allocation entry
1664           * since the secassoc is now unique and
1665           * can't be allocated to any other socket
1666           */
```

26

```
1667         DPRINTF(IDL MAJOR EVENT,("key alloc: making keynode unique..."));
1668         keynode->secassoc->state |= K_UNIQUE;
1669         prevnp->next = np->next;
1670         KFree(np);
1671         DPRINTF(IDL_MAJOR_EVENT,("done\n"));
1672       }
1673     keynode->secassoc->state |= K USED;
1674     keynode->secassoc->state |= K_OUTBOUND;
1675     keynode->alloc_count++;
1676
1677     /*
1678      * Add socket to list of socket using secassoc.
1679      */
1680     DPRINTF(IDL GROSS EVENT,("key_alloc: adding so to solist..."));
1681     newsp->socket = socket;
1682     newsp->next = keynode->solist->next;
1683     keynode->solist->next = newsp;
1684     DPRINTF(IDL_GROSS EVENT,("done\n"));
1685     *keynodep = keynode;
1686     splx(s);
1687     return(0);
1688   }
1689   *keynodep = NULL;
1690   return(0);
1691 }
1692
1693
1694 /*------------------------------------------------------------------
1695  * key free():
1696  *       Decrement the refcount for a key table entry.  If the entry is
1697  *       marked dead, and the refcount is zero, we go ahead and delete
        it.
1698   ---------------------------------------------------------------------*
     /
1699 void
1700 key_free(keynode)
1701       struct key_tblnode *keynode;
1702 {
1703   DPRINTF(IDL MAJOR EVENT,("Entering key_free
     w/keynode=0x%x\n",keynode));
1704   if (!keynode) {
1705     DPRINTF(IDL_CRITICAL,("Warning: key_free got null pointer\n"));
1706     return;
1707   }
1708   (keynode->ref count)--;
1709   if (keynode->ref count < 0) {
1710     DPRINTF(IDL CRITICAL,("Warning: key_free decremented refcount to
     %d\n",keynode->ref_count));
1711   }
1712   if ((keynode->secassoc->state & K_DEAD) && (keynode->ref_count <= 0))
     {
1713     DPRINTF(IDL MAJOR EVENT,("key free: calling key_delete\n"));
1714     key_delete(keynode->secassoc);
1715   }
1716 }
1717
1718 /*---------------------------------------------------------------
1719  * getassocbyspi():
1720  *       Get a security association for a given type, src, dst, and spi.
1721  *
1722  *       Returns: 0 if sucessfull
1723  *                -1 if error/not found
1724  *
1725  *       Caller must convert spi to host order.  Function assumes spi is
1726  *       in host order!
1727   -----------------------------------------------------------------*
```

```
1727   /
1728   int
1729   getassocbyspi(type, src, dst, spi, keyentry)
1730        u int type;
1731        struct sockaddr *src;
1732        struct sockaddr *dst;
1733        u int32 spi;
1734        struct key_tblnode **keyentry;
1735   {
1736     char buf[MAXHASHKEYLEN];
1737     int len, indx;
1738     struct key_tblnode *keynode, *prevkeynode = 0;
1739
1740     DPRINTF(IDL_GROSS_EVENT,("Entering getassocbyspi w/type=%u spi=%u\n",
          type,spi));
1741
1742     /*
1743      *  We treat esp-transport mode and esp-tunnel mode
1744      *  as a single type in the keytable.
1745      */
1746     if (type == SS ENCRYPTION NETWORK)
1747        type = SS_ENCRYPTION_TRANSPORT;
1748
1749     *keyentry = NULL;
1750     bzero(&buf, sizeof(buf));
1751     len = key createkey((char *)&buf, type, src, dst, spi, 0);
1752     indx = key gethashval((char *)&buf, len, KEYTBLSIZE);
1753     DPRINTF(IDL FINISHED,("getassocbyspi: indx=%d\n",indx));
1754     DDO(IDL_FINISHED,dump sockaddr(src);dump sockaddr(dst));
1755     keynode = key search(type, src, dst, spi, indx, &prevkeynode);
1756     DPRINTF(IDL GROSS EVENT,("getassocbyspi: keysearch
          ret=0x%x\n",keynode));
1757     if (keynode && !(keynode->secassoc->state & (K DEAD | K LARVAL))) {
1758        DPRINTF(IDL EVENT,("getassocbyspi: found secassoc!\n"));
1759        (keynode->ref count)++;
1760        *keyentry = keynode;
1761     } else {
1762        DPRINTF(IDL MAJOR_EVENT,("getassocbyspi: secassoc not found!\n"));
1763        return (-1);
1764     }
1765     return(0);
1766   }
1767
1768
1769   /*-------------------------------------------------------------------
1770    * getassocbysocket():
1771    *      Get a security association for a given type, src, dst, and
          socket.
1772    *      If not found, try to allocate one.
1773    *      Returns: 0 if successfull
1774    *              -1 if error condition/secassoc not found (*keyentry =
          NULL)
1775    *               1 if secassoc temporarily unavailable (*keynetry =
          NULL)
1776    *                 (e.g., key mgnt. daemon(s) called)
1777    -----------------------------------------------------------------*
       /
1778   int
1779   getassocbysocket(type, src, dst, socket, unique_key, keyentry)
1780        u int type;
1781        struct sockaddr *src;
1782        struct sockaddr *dst;
1783        struct socket *socket;
1784        u int unique key;
1785        struct key_tblnode **keyentry;
1786   {
1787     struct key_tblnode *keynode = 0;
```

28

```
1788    struct key so2spinode *np;
1789    int len, indx;
1790    u int32 spi;
1791    u_int realtype;
1792
1793    DPRINTF(IDL GROSS EVENT,("Entering getassocbysocket w/type=%u
        so=0x%x\n",type,socket));
1794
1795    /*
1796     *  We treat esp-transport mode and esp-tunnel mode
1797     *  as a single type in the keytable.  This has a side
1798     *  effect that socket using both esp-transport and
1799     *  esp-tunnel will use the same security association
1800     *  for both modes.  Is this a problem?
1801     */
1802    realtype = type;
1803    if (type == SS ENCRYPTION NETWORK)
1804      type = SS_ENCRYPTION_TRANSPORT;
1805
1806    if (np = key sosearch(type, src, dst, socket)) {
1807      if (np->keynode && np->keynode->secassoc &&
1808      !(np->keynode->secassoc->state & (K DEAD | K LARVAL))) {
1809        DPRINTF(IDL FINISHED,("getassocbysocket: found secassoc!\n"));
1810        (np->keynode->ref count)++;
1811        *keyentry = np->keynode;
1812        return(0);
1813      }
1814    }
1815
1816    /*
1817     * No secassoc has been allocated to socket,
1818     * so allocate one, if available
1819     */
1820    DPRINTF(IDL EVENT,("getassocbyso: can't find it, trying to
        allocate!\n"));
1821    if (key_alloc(realtype, src, dst, socket, unique_key, &keynode) == 0)
        {
1822      if (keynode) {
1823        DPRINTF(IDL EVENT,("getassocbyso: key_alloc found secassoc!\n"));
1824        keynode->ref count++;
1825        *keyentry = keynode;
1826        return(0);
1827      } else {
1828        /*
1829         * Kick key mgnt. daemon(s)
1830         * (this should be done in ipsec output policy() instead or
1831         * selectively called based on a flag value)
1832         */
1833        DPRINTF(IDL FINISHED,("getassocbyso: calling key mgnt
            daemons!\n"));
1834        *keyentry = NULL;
1835        if (key acquire(realtype, src, dst) == 0)
1836      return (1);
1837        else
1838      return(-1);
1839      }
1840    }
1841    *keyentry = NULL;
1842    return(-1);
1843  }
1844
1845
```

```
 1   /*------------------------------------------------------------------------
 2    * key.h :      Declarations and Definitions for Key Engine for BSD.
 3    *
 4    * Copyright 1995 by Bao Phan, Randall Atkinson, & Dan McDonald,
 5    * All Rights Reserved.  All rights have been assigned to the US
 6    * Naval Research Laboratory (NRL).  The NRL Copyright Notice and
 7    * License Agreement governs distribution and use of this software.
 8    *
 9    *  Patents are pending on this technology.  NRL grants a license
10    *  to use this technology at no cost under the terms below with
11    * the additional requirement that software, hardware, and
12    * documentation relating to use of this technology must include
13    * the note that:
14    *       This product includes technology developed at and
15    *       licensed from the Information Technology Division,
16    * US Naval Research Laboratory.
17    *
18    ------------------------------------------------------------------------*/
19   /*------------------------------------------------------------------------
20   #   @(#)COPYRIGHT   1.1a (NRL) 17 August 1995
21
22   COPYRIGHT NOTICE
23
24   All of the documentation and software included in this software
25   distribution from the US Naval Research Laboratory (NRL) are
26   copyrighted by their respective developers.
27
28   This software and documentation were developed at NRL by various
29   people.  Those developers have each copyrighted the portions that they
30   developed at NRL and have assigned All Rights for those portions to
31   NRL.  Outside the USA, NRL also has copyright on the software
32   developed at NRL. The affected files all contain specific copyright
33   notices and those notices must be retained in any derived work.
34
35   NRL LICENSE
36
37   NRL grants permission for redistribution and use in source and binary
38   forms, with or without modification, of the software and documentation
39   created at NRL provided that the following conditions are met:
40
41   1. Redistributions of source code must retain the above copyright
42      notice, this list of conditions and the following disclaimer.
43   2. Redistributions in binary form must reproduce the above copyright
44      notice, this list of conditions and the following disclaimer in the
45      documentation and/or other materials provided with the distribution.
46   3. All advertising materials mentioning features or use of this software
47      must display the following acknowledgement:
48
49      This product includes software developed at the Information
50      Technology Division, US Naval Research Laboratory.
51
52   4. Neither the name of the NRL nor the names of its contributors
53      may be used to endorse or promote products derived from this software
54      without specific prior written permission.
55
56   THE SOFTWARE PROVIDED BY NRL IS PROVIDED BY NRL AND CONTRIBUTORS ``AS
57   IS'' AND ANY EXPRESS OR IMPLIED WARRANTIES, INCLUDING, BUT NOT LIMITED
58   TO, THE IMPLIED WARRANTIES OF MERCHANTABILITY AND FITNESS FOR A
59   PARTICULAR PURPOSE ARE DISCLAIMED.  IN NO EVENT SHALL NRL OR
60   CONTRIBUTORS BE LIABLE FOR ANY DIRECT, INDIRECT, INCIDENTAL, SPECIAL,
61   EXEMPLARY, OR CONSEQUENTIAL DAMAGES (INCLUDING, BUT NOT LIMITED TO,
62   PROCUREMENT OF SUBSTITUTE GOODS OR SERVICES; LOSS OF USE, DATA, OR
63   PROFITS; OR BUSINESS INTERRUPTION) HOWEVER CAUSED AND ON ANY THEORY OF
64   LIABILITY, WHETHER IN CONTRACT, STRICT LIABILITY, OR TORT (INCLUDING
65   NEGLIGENCE OR OTHERWISE) ARISING IN ANY WAY OUT OF THE USE OF THIS
66   SOFTWARE, EVEN IF ADVISED OF THE POSSIBILITY OF SUCH DAMAGE.
67
```

1

```
 68   The views and conclusions contained in the software and documentation
 69   are those of the authors and should not be interpreted as representing
 70   official policies, either expressed or implied, of the US Naval
 71   Research Laboratory (NRL).
 72
 73   -------------------------------------------------------------------*/
 74
 75
 76   /*
 77    * PF_KEY messages
 78    */
 79
 80   #define KEY_ADD            1
 81   #define KEY_DELETE         2
 82   #define KEY_UPDATE         3
 83   #define KEY_GET            4
 84   #define KEY_ACQUIRE        5
 85   #define KEY_GETSPI         6
 86   #define KEY_REGISTER       7
 87   #define KEY_EXPIRE         8
 88   #define KEY_DUMP           9
 89   #define KEY_FLUSH          10
 90
 91   #define KEY_VERSION        1
 92   #define POLICY_VERSION     1
 93
 94   /*
 95    * Security association state
 96    */
 97
 98   #define K_USED             0x1     /* Key used/not used */
 99   #define K_UNIQUE           0x2     /* Key unique/reusable */
100   #define K_LARVAL           0x4     /* SPI assigned, but sa incomplete */
101   #define K_ZOMBIE           0x8     /* sa expired but still useable */
102   #define K_DEAD             0x10    /* sa marked for deletion, ready for
      reaping */
103   #define K_INBOUND          0x20    /* sa for inbound packets, ie. dst=myhost
      */
104   #define K_OUTBOUND         0x40    /* sa for outbound packets, ie.
      src=myhost */
105
106   /*
107    * Structure for key message header.
108    * PF_KEY message consists of key_msghdr followed by
109    * src sockaddr, dest sockaddr, from sockaddr, key, and iv.
110    * Assumes size of key message header less than MHLEN.
111    */
112
113   struct key_msghdr {
114      u_short key_msglen;     /* length of message including
      src/dst/from/key/iv */
115      u_char  key_msgvers;    /* key version number */
116      u_char  key_msgtype;    /* key message type, eg. KEY_ADD */
117      pid_t   key_pid;        /* process id of message sender */
118      int     key_seq;        /* message sequence number */
119      int     key_errno;      /* error code */
120      u_int8  type;           /* type of security association */
121      u_int8  state;          /* state of security association */
122      u_int8  label;          /* sensitivity level */
123      u_int32 spi;            /* spi value */
124      u_int8  keylen;         /* key length */
125      u_int8  ivlen;          /* iv length */
126      u_int8  algorithm;      /* algorithm identifier */
127      u_int8  lifetype;       /* type of lifetime */
128      u_int32 lifetime1;      /* lifetime value 1 */
129      u_int32 lifetime2;      /* lifetime value 2 */
130   };
```

2

```
131
132   struct key msgdata {
133     struct sockaddr *src;      /* source host address */
134     struct sockaddr *dst;      /* destination host address */
135     struct sockaddr *from;     /* originator of security association */
136     caddr t iv;                /* initialization vector */
137     caddr t key;               /* key */
138     int ivlen;                 /* key length */
139     int keylen;                /* iv length */
140   };
141
142   struct policy msghdr {
143     u short policy msglen;     /* message length */
144     u char  policy msgvers;    /* message version */
145     u char  policy msgtype;    /* message type */
146     int     policy seq;        /* message sequence number */
147     int     policy_errno;      /* error code */
148   };
149
150
151   #ifdef KERNEL
152
153   /*
154    * Key engine table structures
155    */
156
157   struct socketlist {
158     struct socket *socket;       /* pointer to socket */
159     struct socketlist *next;     /* next */
160   };
161
162   struct key tblnode {
163     int alloc count;             /* number of sockets allocated to
          secassoc */
164     int ref_count;               /* number of sockets referencing secassoc
          */
165     struct socketlist *solist;   /* list of sockets allocated to secassoc
          */
166     struct ipsec assoc *secassoc; /* security association */
167     struct key_tblnode *next;    /* next node */
168   };
169
170   struct key allocnode {
171     struct key tblnode *keynode;
172     struct key_allocnode *next;
173   };
174
175   struct key so2spinode {
176     struct socket *socket;       /* socket pointer */
177     struct key_tblnode *keynode; /* pointer to tblnode containing secassoc
          */
178                   /*  info for socket  */
179     struct key_so2spinode *next;
180   };
181
182   struct key registry {
183     u int8 type;               /* secassoc type that key mgnt. daemon can
          acquire */
184     struct socket *socket;     /* key management daemon socket pointer */
185     struct key_registry *next;
186   };
187
188   struct key acquirelist {
189     u int8 type;                 /* secassoc type to acquire */
190     struct sockaddr_in6 target;  /* destination address of secassoc */
191     u int32 count;               /* number of acquire messages sent */
192     u_long expiretime;           /* expiration time for acquire message */
```

3

```
193      struct key_acquirelist *next;
194   };
195
196   struct keyso_cb {
197      int ip4_count;           /* IPv4 */
198      int ip6_count;           /* IPv6 */
199      int any_count;           /* Sum of above counters */
200   };
201
202   #endif
203
204   /*
205    * Useful macros
206    */
207
208   #ifndef KERNEL
209   #define K_Malloc(p, t, n) (p = (t) malloc((unsigned int)(n)))
210   #define KFree(p) free((char *)p);
211   #else
212   #define K_Malloc(p, t, n) (p = (t) malloc((unsigned long)(n), M_SECA,
      M_DONTWAIT))
213   #define KFree(p) free((caddr_t)p, M_SECA);
214   #endif /* KERNEL */
215
216   #ifdef KERNEL
217   void    key_init    _P((void));
218   void    key_cbinit   _P((void));
219   void    key_inittables    _P((void));
220   int     key_secassoc2msghdr __P((struct ipsec_assoc *, struct key_msghdr
      *,
221                   struct key_msgdata *));
222   int     key_msghdr2secassoc __P((struct ipsec_assoc *, struct key_msghdr
      *,
223                   struct key_msgdata *));
224   int     key_add    P((struct ipsec_assoc *));
225   int     key_delete   P((struct ipsec_assoc *));
226   int     key_get    P((u_int, struct sockaddr *, struct sockaddr *, u_int32,
227               struct ipsec_assoc **));
228   void    key_flush   P((void));
229   int     key_dump   _P((struct socket *));
230   int     key_getspi   P((u_int, struct sockaddr *, struct sockaddr *,
231                   u_int32 *));
232   int     key_update   _P((struct ipsec_assoc *));
233   int     key_register   _P((struct socket *, u_int));
234   void    key_unregister   P((struct socket *, u_int, int));
235   int     key_acquire   _P((u_int, struct sockaddr *, struct sockaddr *));
236   int     getassocbyspi   P((u_int, struct sockaddr *, struct sockaddr *,
237                   u_int32, struct key_tblnode **));
238   int     getassocbysocket   P((u_int, struct sockaddr *, struct sockaddr *,
239                   struct socket *, u_int, struct key_tblnode **));
240   void    key_free   _P((struct key_tblnode *));
241   int     key_output   P((struct mbuf *, struct socket *));
242   int     key_usrreq   __P((struct socket *, int, struct mbuf *, struct mbuf
      *,
243                   struct mbuf *));
244   #endif
245
```

4