Sun RPC Information

Appendix 1

A1.1 RPC Source Code Availability

The source code for the SunOS 4.0 RPC and XDR Library, referred to as RPCSRC 4.0, is a license-free version of Sun's RPC and XDR Library. RPCSRC is available free to anyone with access to the Internet via an anonymous file-transfer program (ftp) login from one of the RPCSRC archive sites. Because the archive sites change over time, consult your local system administrator for more information on locating an archive site and on using the Internet. You can also obtain RPCSRC 4.0 directly from Sun for a nominal processing fee. The part numbers are:

RPC-4.0-X-X-5 RPCSRC on 1/4-inch tape RPC-4.0-X-X-6 RPCSRC on 1/2-inch tape

In the future, new releases of the RPC and XDR Library will be available from RPCSRC archive sites.

A1.2 RPC Program Numbers

Program Number

Program numbers are assigned in groups of 0x20000000 according to the following chart:

Description

| | - coorpaon |
|-------------------------|-----------------------------|
| 0x00000000 - 0x1FFFFFFF | Defined by Sun Microsystems |
| 0x20000000 - 0x3FFFFFFF | Defined by User |
| 0x40000000 - 0x5FFFFFFF | Transient |
| 0x60000000 - 0x7FFFFFFF | Reserved |
|)x80000000 - 0x9FFFFFFF | Reserved |
|)xA0000000 - 0xBFFFFFFF | Reserved |
|)xC0000000 - 0xDFFFFFFF | Reserved |
| 0xE0000000 - 0xFFFFFFFF | Reserved |

Sun Microsystems administers the first group of numbers, which should be identical for all users of Sun's RPC Library, to ensure that the global program numbers are unique. If you develop an application that might be of general interest, or that might become a product, then you should obtain a program number from the first range. The second group of numbers is reserved for applications specific to you, in that these applications are only running on your network. This range is intended primarily for developing new programs. If you develop internal distributed applications, then somebody within your organization should maintain these numbers to ensure that two different applications do not try to use the same program number. The third group is reserved for applications that generate program numbers dynamically. such as applications that use Callback RPC. The final groups are reserved for future use, and should not be used. Blocks of numbers are also available for assignment to companies for use internally or for assignment to your customers.

To obtain a unique RPC program number and to optionally register a protocol specification, send a request by electronic mail to *rpc@sun* or write to: RPC Administrator, Sun Microsystems, 2550 Garcia Ave., Mountain View, CA 94043

Below is a list of the RPC program numbers that have been assigned and are public. On Unix systems, these program numbers are usually found in the file */etc/rpc*.

| portmapper | 100000 portmap sunrpc |
|---------------|--------------------------------|
| rstatd | 100001 rstat rup perfmeter |
| rusersd | 100002 rusers |
| nfs | 100003 nfsprog |
| ypserv | 100004 ypprog |
| mountd | 100005 mount showmount |
| ypbind | 100007 |
| walld | 100008 rwall shutdown |
| yppasswdd | 100009 yppasswd |
| etherstatd | 100010 etherstat |
| rquotad | 100011 rquotaprog quota rquota |
| sprayd | 100012 spray |
| 3270_mapper | 100013 |
| rje_mapper | 100014 |
| selection_svc | 100015 selnsvc |
| database_svc | 100016 |
| rexd | 100017 rex |
| alis | 100018 |

| sched | 100019 |
|------------|--------------------|
| llockmgr | 100020 |
| nlockmgr | 100021 |
| x25.inr | 100022 |
| statmon | 100023 |
| status | 100024 |
| bootparam | 100026 |
| ypupdated | 100028 ypupdate |
| keyserv | 100029 keyserver |
| tfsd | 100037 |
| nsed | 100038 |
| nsemntd | 100039 |
| Netlicense | 100062 rpc.netlicd |
| | - |

The following chart lists the currently used authentication numbers:

| Authentication Number | Description |
|-----------------------|-----------------------|
| 0 | None |
| 1 | UNIX-style |
| 2 | Short hand UNIX-style |
| 3 | DES |

Sun Microsystems administers the entire range of authentication numbers. If you develop a new authentication flavor and wish to reserve the authentication number, then you should obtain a unique authentication number. Blocks of numbers are also available for assignment to companies. The procedure for obtaining an authentication number is the same as for obtaining a program number. External Data Representation Standard: Protocol Specification

A2.1 Status of This Standard

Note: This appendix specifies a protocol that Sun Microsystems, Inc., and others are using. It has been designated RFC1014 by the ARPA Network Information Center. **Appendix**

2

A2.2 Introduction

XDR is a standard for the description and encoding of data. It is useful for transferring data between different computer architectures, and has been used to communicate data between such diverse machines as the Sun Workstation, VAX, IBM-PC, and Cray. XDR fits into the ISO presentation layer, and is roughly analogous in purpose to X.409, ISO Abstract Syntax Notation. The major difference between these two is that XDR uses implicit typing, while X.409 uses explicit typing.

XDR uses a language to describe data formats. The language can be used only to describe data; it is not a programming language. This language allows one to describe intricate data formats in a concise manner. The alternative of using graphical representations (itself an informal language) quickly becomes incomprehensible when faced with complexity. The XDR language itself is similar to the C language [1], just as Courier [4] is similar to Mesa. Protocols such as Sun RPC (Remote Procedure Call) and the NFS (Network File System) use XDR to describe the format of their data.

The XDR standard makes the following assumption: that bytes (or octets) are portable, where a byte is defined to be 8 bits of data. A given hardware device should encode the bytes onto the various media in such a way that other hardware devices may decode the bytes without loss of meaning. For example, the Ethernet standard suggests that bytes be encoded in "little-endian" style [2], or least significant bit first.

Basic Block Size

The representation of all items requires a multiple of four bytes (or 32 bits) of data. The bytes are numbered 0 through n - 1. The bytes are read or written to some byte stream such that byte m always precedes byte m + 1. If the n bytes needed to contain the data are not a multiple of four, then the n bytes are followed by enough (0 to 3) residual zero bytes, r, to make the total byte count a multiple of 4.

We include the familiar graphic box notation for illustration and comparison. In most illustrations, each box (delimited by a plus sign at the 4 corners and vertical bars and dashes) depicts a byte. Ellipses (...) between boxes show zero or more additional bytes where required.

A Block

| ++- byte 0 | byte 1 | · · · · + · · · by | +- te n-1 | 0 | -+ | ++ |
|-----------------|--------|-------------------------|----------------|-------|-----------|----|
| < | n byte | s r (where (| > < n + r mod | 4 = 0 | - r bytes | 5> |

A2.3 XDR Data Types

Each of the sections that follow describes a data type defined in the XDR standard, shows how it is declared in the language, and includes a graphic illustration of its encoding.

For each data type in the language we show a general paradigm declaration. Note that angle brackets ((and)) denote variable length sequences of data and square brackets ([and]) denote fixedlength sequences of data. "n", "m" and "r" denote integers. For the full language specification and more formal definitions of terms such as "identifier" and "declaration", refer to The XDR Language Specification presented later.

For some data types, more specific examples are included. A more extensive example of a data description is in *An Example of an XDR Data Description* that follows.

Integer

An XDR signed integer is a 32-bit datum that encodes an integer in the range [-2147483648,2147483647]. The integer is represented in two's complement notation. The most and least significant bytes are 0 and 3, respectively. Integers are declared as follows:

Integer

| (MSB) | | (LSB) | | | | |
|--------|--------|--------|--------|---|--|--|
| byte 0 | byte 1 | byte 2 | byte 3 | | | |
| < | 32 | bits | : | > | | |

Unsigned Integer

An XDR unsigned integer is a 32-bit datum that encodes a nonnegative integer in the range [0,4294967295]. It is represented by an unsigned binary number with most and least significant bytes of 0 and 3, respectively. An unsigned integer is declared as follows:

Unsigned Integer

| (MSB) | | (LSB) | | | | |
|--------|--------|--------|--------|--|--|--|
| byte 0 | byte 1 | byte 2 | byte 3 | | | |
| < | 32 | bits | > | | | |

Enumeration

Enumerations have the same representation as signed integers. Enumerations are handy for describing subsets of the integers. Enumerated data is declared as follows:

enum { name-identifier = constant, ... } identifier;

For example, the three colors red, yellow, and blue could be described by an enumerated type:

enum { RED = 2, YELLOW = 3, BLUE = 5 } colors;

It is an error to encode as an enum any other integer than those that have been given assignments in the enum declaration.

Boolean

Booleans are important enough and occur frequently enough to warrant their own explicit type in the standard. Booleans are declared as follows:

bool identifier;

This is equivalent to:

enum { FALSE = 0, TRUE = 1 } identifier;

Hyper Integer and Unsigned Hyper Integer

The standard also defines 64-bit (8-byte) numbers called hyper integer and unsigned hyper integer. Their representations are the obvious extensions of integer and unsigned integer defined above. They are represented in two's complement notation. The most and least significant bytes are 0 and 7, respectively. Their declarations:

Hyper Integer Unsigned Hyper Integer

| (MSB) | (LSB) |
|---|-----------------|
| byte 0 byte 1 byte 2 byte 3 byte 4 byte 5 | byte 6 byte 7 |
| <64 bits | > |

Floating-point

The standard defines the floating-point data type "float" (32 bits or 4 bytes). The encoding used is the IEEE standard for normalized single-precision floating-point numbers [3]. Three fields describe the single-precision floating-point number. The first field, which takes ((uses?)) one bit, refers to the sign of the number; values 0 and 1 represent positive and negative, respectively. The second field describes the exponent of the number to

the base 2, with the exponent biased by 127. Eight bits are devoted to the second field. The third field refers to the fractional part of the number's mantissa with base 2; 23 bits are devoted to this field.

Therefore, the floating-point number is described by:

(-1)**S * 2**(E-Bias) * 1.F

It is declared as follows:

Single-Precision Floating-Point

| + - S | byte 0 E | byte 1 byte 2 byte 3 F |
|---------------|-----------------|---------------------------------|
| 1 < | <- 8 -> | < 23 bits > 32 bits > |

Just as the most and least significant bytes of a number are 0 and 3, the most and least significant bits of a single-precision floating-point number are 0 and 31. The beginning bit (and most significant bit) offsets of S, E, and F are 0, 1, and 9, respectively. Note that these numbers refer to the mathematical positions of the bits, and NOT to their actual physical locations (which vary from medium to medium).

The IEEE specifications should be consulted concerning the encoding for signed zero, signed infinity (overflow), and denormalized numbers (underflow) [3]. According to IEEE specifications, the "NaN" (not a number) is system dependent and should not be used externally.

Double-precision Floating-point

The standard defines the encoding for the double-precision floating-point data type "double" (64 bits or 8 bytes). The encoding used is the IEEE standard for normalized double-precision floating-point numbers [3]. The standard encodes the following three fields, which describe the double-precision floating-point number. The first field, which takes one bit, refers to the sign of the number; values 0 and 1 represent positive and negative, respectively. The second field describes the exponent of the number to the base 2, with the exponent biased by 1023. Eleven bits are devoted to the second field. The third field refers to the fractional part of the number's mantissa with base 2; 52 bits are devoted to this field. Therefore, the floating-point number is described by:

(-1)**S * 2**(E-Bias) * 1.F

It is declared as follows:

Double-Precision Floating-Point

| + - | byte 0 | byte 1 byte 2 byte 3 byte 4 byte 5 byte 6 byte 7 |
|-----|----------|--|
| Ś | Е | F |
| 1 | <- 11 -> | <> 52 bits> |
| < - | | 64 bits > |

Just as the most and least significant bytes of a number are 0 and 3, the most and least significant bits of a double-precision floating-point number are 0 and 63. The beginning bit (and most significant bit) offsets of S, E, and F are 0, 1, and 12, respectively. Note that these numbers refer to the mathematical positions of the bits, and NOT to their actual physical locations (which vary from medium to medium).

The IEEE specifications should be consulted concerning the encoding for signed zero, signed infinity (overflow), and denormalized numbers (underflow) [3]. According to IEEE specifications, the "NaN" (not a number) is system dependent and should not be used externally.

Fixed-length Opaque Data

At times, fixed-length uninterpreted data needs to be passed among machines. This data is called "opaque" and is declared as follows:

```
opaque identifier[n];
```

where the constant n is the (static) number of bytes necessary to contain the opaque data. If n is not a multiple of four, then the n bytes are followed by enough (0 to 3) residual zero bytes, r, to make the total byte count of the opaque object a multiple of four.

Fixed-Length Opaque

| 0 | | 1 | ••• | | ± | L | | 1 | |
|---|--------|------|---------------|---------|---------------|----------|-----------|---|---|
| | byte 0 | byte | 1 | ••• | byte n – 1 | 0 | | 0 | |
| ŀ | < | n | bytes | · · · · | > | < | - r bytes | > | + |
| · | < | | $\cdot n + r$ | (whe | ere (n + r) m | od 4 = 0 |) | > | |

Variable-length Opaque Data

The standard also provides for variable-length (counted) opaque data, defined as a sequence of n (numbered 0 through n - 1) arbitrary bytes to be the number n encoded as an unsigned integer (as described below), and followed by the n bytes of the sequence.

Byte m of the sequence always precedes byte m + 1 of the sequence, and byte 0 of the sequence always follows the sequence's length (count). Enough (0 to 3) residual zero bytes, r, to make the total byte count a multiple of four. Variable-length opaque data is declared in the following way:

opaque identifier $\langle m \rangle$;

or

opaque identifier $\langle \rangle$;

The constant m denotes an upper bound of the number of bytes that the sequence may contain. If m is not specified, as in the second declaration, it is assumed to be $(2^{**}32) - 1$, the maximum length. The constant m would normally be found in a protocol specification. For example, a filing protocol may state that the maximum data transfer size is 8192 bytes, as follows:

opaque filedata $\langle 8192 \rangle$;

This can be illustrated as follows:

Variable-Length Opaque

| 0 | 1 | 2 | 3 | 4 | 5 | • • • | | | | | | |
|----|---------|---------|--------|------------|----------|--------|---------|----------|---|---------|---|---|
| + | + le | ength 1 | + n | + byt | e 0 b | yte 1 | + + | n-1 | 0 | + + | 0 | + |
| <- | 4 | 4 bytes | ; | -> < | +- | -n byt | es | · > · | < | r bytes | > | ļ |
| <- | | | 4 | +n+r | (where | (n+r) |) mod | 4 = 0) - | | | > | |

It is an error to encode a length greater than the maximum described in the specification.

String

The standard defines a string of n (numbered 0 through n - 1) ASCII bytes to be the number n encoded as an unsigned integer (as described above), and followed by the n bytes of the string. Byte m of the string always precedes byte m + 1 of the string, and byte 0 of the string always follows the string's length. If n is not a multiple of four, then the n bytes are followed by enough (0 to 3) residual zero bytes, r, to make the total byte count a multiple of four. Counted byte strings are declared as follows:

string object(m);

or

string object();

The constant m denotes an upper bound of the number of bytes that a string may contain. If m is not specified, as in the second declaration, it is assumed to be $(2^{**}32) - 1$, the maximum length. The constant m would normally be found in a protocol specification. For example, a filing protocol may state that a file name can be no longer than 255 bytes, as follows:

string filename(255);

Which can be illustrated as:

A String

| 0 | 1 | 2 | 3 | 4 | 5 | | | | | | | 4 |
|----|----|--------|--------|------|----------|--------|-------|--------|---|---|---|---|
| +- | le | ngth r | + 1 | by | te 0 b | yte 1 | | n-1 | 0 | $\begin{vmatrix} & \cdot & \cdot & \cdot \\ & \cdot & \cdot & \cdot \\ & \cdot & \cdot & \cdot$ | 0 | |
| < | 4 | bytes | | -> < | | -n byt | es | > | < | r bytes | > | 1 |
| < | | | 4 | +n+r | (where | (n+r) |) mod | 4 = 0) | | | > | |

It is an error to encode a length greater than the maximum described in the specification.

Fixed-length Array

Declarations for fixed-length arrays of homogeneous elements are in the following form:

type-name identifier[n];

Fixed-length arrays of elements numbered 0 through n - 1 are encoded by individually encoding the elements of the array in their natural order, 0 through n - 1. Each element's size is a multiple of four bytes. Though all elements are of the same type, the elements may have different sizes. For example, in a fixedlength array of strings, all elements are of type "string", yet each element will vary in its length.

Fixed-Length Array

| ++++ | ++++ | + + + + + | + |
|-----------|------------|--------------------------|---|
| element 0 | element 1 | $ \ldots $ element $n-1$ | 1 |
| +++ | + + + + | + +++ | + |
| < | n elements | | > |

Variable-length Array

Counted arrays provide the ability to encode variable-length arrays of homogeneous elements. The array is encoded as the element count n (an unsigned integer) followed by the encoding of each of the array's elements, starting with element 0 and progressing through element n - 1. The declaration for variable-length arrays follows this form:

type-name identifier(m);

or

type-name identifier $\langle \rangle$;

The constant m specifies the maximum acceptable element count of an array; if m is not specified, as in the second declaration, it is assumed to be $(2^{**}32) - 1$.

Counted Array

| 0 | 1 | 2 | 3 | | | | |
|-----|-------|-------|-----|---------------|------------|-----------|---------------|
| + - | - + - | - + - | -+ | ++++ | ++ | + • • • - | + + + + |
| | | n | | element 0 | element 1 | | element n – 1 |
| + - | - + - | - + - | -+ | + ~ _ + + + + | ++ | + | ·+++ |
| < | 4 | byte | es> | < | n elements | s | > |

It is an error to encode a value of n that is greater than the maximum described in the specification.

Structure

Structures are declared as follows:

struct {

```
component-declaration-A; component-declaration-B;
```

} identifier;

. . .

The components of the structure are encoded in the order of their declaration in the structure. Each component's size is a multiple of four bytes, though the components may be different sizes.

Structure

| + | | • | • |
|---|---|---|---|
| component A component B | | | • |
| + + + + + + | ٠ | ٠ | • |

Discriminated Union

A discriminated union is a type composed of a discriminant followed by a type selected from a set of prearranged types according to the value of the discriminant. The type of discriminant is either "int", "unsigned int", or an enumerated type, such as "bool". The component types are called "arms" of the union, and are preceded by the value of the discriminant which implies their encoding. Discriminated unions are declared as follows:

```
union switch (discriminant-declaration) {
    case discriminant-value-A:
    arm-declaration-A;
    case discriminant-value-B:
    arm-declaration-B;
    ...
    default: default-declaration;
} identifier;
```

Each "case" keyword is followed by a legal value of the discriminant. The default arm is optional. If the discriminant is not specified, then a valid encoding of the union cannot take on unspecified discriminant values. The size of the implied arm is always a multiple of four bytes.

The discriminated union is encoded as its discriminant followed by the encoding of the implied arm.

```
Discriminated Union
```

| 0 | 1 | 2 | 3 | |
|-----|-----|---------|-----|-------------|
| | dis | crimin | ant | implied arm |
| < - | | 4 bytes | s | > |

Void

An XDR void is a 0-byte quantity. Voids are useful for describing operations that take no data as input or no data as output. They are also useful in unions, where some arms may contain data and others do not. The declaration is simply as follows:

void;

Voids are illustrated as follows:

```
+++
||
+++
--><-- 0 bytes
```

Constant

The data declaration for a constant follows this form:

const name-identifier = n;

"const" is used to define a symbolic name for a constant; it does not declare any data. The symbolic constant may be used anywhere a regular constant may be used. For example, the following defines a symbolic constant DOZEN, equal to 12.

```
const DOZEN = 12;
```

Typedef

"typedef" does not declare any data either, but serves to define new identifiers for declaring data. The syntax is:

typedef declaration;

The new type name is actually the variable name in the declaration part of the typedef. For example, the following defines a new type called "eggbox" using an existing type called "egg":

```
typedef egg eggbox[DOZEN];
```

Variables declared using the new type name have the same type as the new type name would have in the typedef, if it were considered a variable. For example, the following two declarations are equivalent in declaring the variable "fresheggs":

```
eggbox fresheggs;
egg fresheggs[DOZEN];
```

When a typedef involves a struct, enum, or union definition, there is another (preferred) syntax that may be used to define the same type. In general, a typedef of the following form:

typedef ((struct, union, or enum definition)) identifier;

may be converted to the alternative form by removing the "typedef" part and placing the identifier after the "struct", "union", or "enum" keyword, instead of at the end. For example, here are the two ways to define the type "bool":

```
typedef enum { /* using typedef */
FALSE = 0,
TRUE = 1
} bool;
enum bool { /* preferred alternative */
FALSE = 0,
TRUE = 1
};
```

The reason for preferring this syntax is that one does not have to wait until the end of a declaration to figure out the name of the new type.

Optional-data

Optional-data is one kind of union that occurs so frequently that we give it a special syntax of its own for declaring it. It is declared as follows:

type-name *identifier;

This is equivalent to the following union:

```
union switch (bool opted) {
    case TRUE:
    type-name element;
    case FALSE:
    void;
} identifier;
```

It is also equivalent to the following variable-length array declaration, since the boolean "opted" can be interpreted as the length of the array:

type-name identifier(1);

Optional-data is not so interesting in itself, but it is very useful for describing recursive data-structures such as linked-lists and trees. For example, the following defines a type "stringlist" that encodes lists of arbitrary length strings:

```
struct *stringlist {
        string item();
        stringlist next;
};
```

It could have been equivalently declared as the following union:

```
union stringlist switch (bool opted) {
    case TRUE:
        struct {
            string item{};
            stringlist next;
        } element;
    case FALSE:
        void;
    };
or as a variable-length array:
    struct stringlist(1) {
        string item{};
        stringlist next;
        };
```

Both of these declarations obscure the intention of the stringlist type, so the optional-data declaration is preferred over both of them. The optional-data type also has a close correlation to how recursive data structures are represented in high-level languages such as Pascal or C by use of pointers. In fact, the syntax is the same as that of the C language for pointers.

Areas for Future Enhancement

The XDR standard lacks representations for bit fields and bitmaps, since the standard is based on bytes. Also missing are packed (or binary-coded) decimals. The intent of the XDR standard was not to describe every kind of data that people have ever sent or will ever want to send from machine to machine. Rather, it only describes the most commonly used data-types of high-level languages such as Pascal or C so that applications written in these languages will be able to communicate easily over some medium.

One could imagine extensions to XDR that would let it describe almost any existing protocol, such as TCP. The minimum necessary for this are support for different block sizes and byteorders. The XDR discussed here could then be considered the 4byte big-endian member of a larger XDR family.

A2.4 Discussion

Why a Language for Describing Data?

There are many advantages in using a data-description language such as XDR versus using diagrams. Languages are more formal than diagrams and lead to less ambiguous descriptions of data. Languages are also easier to understand and allow one to think of other issues instead of the low-level details of bit-encoding. Also, there is a close analogy between the types of XDR and a high-level language such as C or Pascal. This makes the implementation of XDR encoding and decoding modules an easier task. Finally, the language specification itself is an ASCII string that can be passed from machine to machine to perform on-the-fly data interpretation.

Why Only One Byte-Order for an XDR Unit?

Supporting two byte-orderings requires a higher level protocol for determining in which byte-order the data is encoded. Since XDR is not a protocol, this cannot be done. However, an advantage of supporting two byte-orderings is that data in XDR format can be written to a magnetic tape, and any machine will be able to interpret it, since no higher level protocol is necessary for determining the byte-order.

Why Does XDR Use Big-Endian Byte-Order?

Yes, it is unfair, but having only one byte-order means you have to be unfair to somebody. Many architectures, such as the Motorola 68000 and IBM 370, support the big-endian byte-order.

Why Is the XDR Unit Four Bytes Wide?

There is a tradeoff in choosing the XDR unit size. Choosing a small size such as two makes the encoded data small, but causes alignment problems for machines that aren't aligned on these boundaries. A large size such as eight means the data will be aligned on virtually every machine, but causes the encoded data to grow too big. We chose four as a compromise. Four is big enough to support most architectures efficiently, except for rare machines such as the eight-byte aligned Cray. Four is also small enough to keep the encoded data restricted to a reasonable size.

Why Must Variable-Length Data Be Padded with Zeros?

It is desirable that the same data encode into the same thing on all machines, so that encoded data can be meaningfully compared or checksummed. Forcing the padded bytes to be zero ensures this.

Why Is There No Explicit Data-Typing?

Data-typing has a relatively high cost for what small advantages it may have. One cost is the expansion of data due to the inserted type fields. Another is the added cost of interpreting these type fields and acting accordingly. Since most protocols already know what type they expect, data-typing supplies only redundant information. However, one can still get the benefits of data-typing using XDR. One way is to encode two things: first a string which is the XDR data description of the encoded data, and then the encoded data itself. Another way is to assign a value to all the types in XDR, and then define a universal type which takes this value as its discriminant and for each value describes the corresponding data type.

A2.5 The XDR Language Specification

Notational Conventions

This specification uses an extended Backus-Naur Form notation for describing the XDR language. Here is a brief description of the notation:

- The characters |, (,), [,], , and * are special.
- Terminal symbols are strings of any characters surrounded by double quotes.
- Non-terminal symbols are strings of non-special characters.
- Alternative items are separated by a vertical bar ("|").
- Optional items are enclosed in brackets.
- Items are grouped together by enclosing them in parentheses.
- \cdot A * following an item means 0 or more occurrences of that item.

For example, consider the following pattern:

"a " "very" (", " " very")* [" cold" "and"] " rainy " ("day" | "night")

An infinite number of strings match this pattern. A few of them are:

"a very rainy day" "a very, very rainy day" "a very cold and rainy day" "a very, very, very cold and rainy night"

Lexical Notes

- Comments begin with '/*' and terminate with '*/'.
- White space serves to separate items and is otherwise ignored.
- An identifier is a letter followed by an optional sequence of letters, digits or underbar ('_'). The case of identifiers is not ignored.

• A constant is a sequence of one or more decimal digits, optionally preceded by a minus-sign ('-').

Syntax Information

```
declaration:
```

```
type-specifier identifier
| type-specifier identifier "[" value "]"
| type-specifier identifier "{" [ value ] "}"
| "opaque" identifier "[" value "]"
| "opaque" identifier "{" [ value ] "}"
| "string" identifier "{" [ value ] "}"
| type-specifier "*" identifier
| "void"
```

value:

```
constant
| identifier
```

type-specifier:

```
[ "unsigned" ] "int"

[ "unsigned" ] "hyper"

"float"

"double"

"bool"

| enum-type-spec

| struct-type-spec

| union-type-spec

| identifier
```

```
enum-type-spec:
```

```
"enum" enum-body
```

```
enum-body
"{"
```

```
( identifier "=" value )
( "," identifier "=" value )*
"}"
```

```
struct-type-spec:
```

```
"struct" struct-body
```

```
struct-body:

"{"

( declaration ";" )

( declaration ";" )*

"}"
```

```
union-type-spec:
         "union" union-body
union-body:
         "switch" "(" declaration ")" "{"
         ( "case" value ":" declaration ";" )
         ("case" value ":" declaration ";")*
         [ "default" ":" declaration ";" ]
         ''}"
constant-def:
         "const" identifier "=" constant ";"
type-def:
         "typedef" declaration ";"
         | "enum" identifier enum-body ";"
| "struct" identifier struct-body ";"
         "'union" identifier union-body
definition:
         type-def
         | constant-def
```

specification: definition *

Syntax Notes

- 1. The following are keywords and cannot be used as identifiers: "bool", "case", "const", "default", "double", "enum", "float", "hyper", opaque", "string", "struct", "switch", "typedef", "union", "unsigned" and "void".
- 2. Only unsigned constants may be used as size specifications for arrays. If an identifier is used, it must have been declared previously as an unsigned constant in a "const" definition.
- 3. Constant and type identifiers within the scope of a specification are in the same name space and must be declared uniquely within this scope.
- 4. Similarly, variable names must be unique within the scope of struct and union declarations. Nested struct and union declarations create new scopes.

5. The discriminant of a union must be of a type that evaluates to an integer. That is, "int", "unsigned int", "bool", an enumerated type or any typedefed type that evaluates to one of these is legal. Also, the case values must be one of the legal values of the discriminant. Finally, a case value may not be specified more than once within the scope of a union declaration.

A2.6 An Example of an XDR Data Description

Here is a short XDR data description of a thing called a "file", which might be used to transfer files from one machine to another.

const MAXUSERNAME = 32; /* max length of a user name */ const MAXFILELEN = 65535; /* max length of a file */ const MAXNAMELEN = 255; /* max length of a file name */ /* * Types of files: enum filekind { /* ascii data */ TEXT = 0, DATA = 1, /* raw data */ EXEC = 2/* executable */ }; /* * File information, per kind of file: */ union filetype switch (filekind kind) { TEXT: case void: /* no extra information */ case DATA: string creator(MAXNAMELEN); /* data creator */ EXEC: case string interpretor(MAXNAMELEN); /* program interpretor */ }; /* * A complete file: */

struct file {

| string filename(MAXNAMELEN); | /* name of file */ |
|------------------------------|-----------------------|
| filetype type; | /* info about file */ |
| string owner (MAXUSERNAME); | /* owner of file */ |
| opaque data(MAXFILELEN); | /* file data */ |
| | |

};

Suppose now that there is a user named "john" who wants to store his lisp program "sillyprog" that contains just the data "(quit)". His file would be encoded as follows:

| Offset | Hex Bytes | ASCII | Description |
|--------|------------------|-------|-----------------------------|
| 0 | 00 00 00 09 | | Length of filename = 9 |
| 4 | 73 69 6c 6c | sill | Filename characters |
| 8 | 79 70 72 6f | ypro | and more characters |
| 12 | $67\ 00\ 00\ 00$ | g | and 3 zero-bytes of fill |
| 16 | 00 00 00 02 | | Filekind is EXEC = 2 |
| 20 | 00 00 00 04 | | Length of interpretor $= 4$ |
| 24 | 6c 69 73 70 | lisp | Interpretor characters |
| 28 | 00 00 00 04 | | Length of owner $= 4$ |
| 32 | 6a 6f 68 6e | john | Owner characters |
| 36 | 00 00 00 06 | | Length of file data = 6 |
| 40 | $28\ 71\ 75\ 69$ | (qui | File data bytes |
| 44 | 74 29 00 00 | t) | and 2 zero-bytes of fill |

References

[1] Brian W. Kernighan & Dennis M. Ritchie, "The C Programming Language", Bell Laboratories, Murray Hill, New Jersey, 1978.

[2] Danny Cohen, "On Holy Wars and a Plea for Peace", IEEE Computer, October 1981.

[3] "IEEE Standard for Binary Floating-Point Arithmetic", ANSI/IEEE Standard 754-1985, Institute of Electrical and Electronics Engineers, August 1985.

[4] "Courier: The Remote Procedure Call Protocol", XEROX Corporation, XSIS 038112, December 1981.

Remote Procedure Call: Protocol Specification

Appendix 3

A3.1 Status of this Memo

Note: This appendix specifies a protocol that Sun Microsystems, Inc., and others are using. It has been designated RFC1057 by the ARPA Network Information Center.

A3.2 Introduction

This appendix specifies a message protocol used in implementing Sun's Remote Procedure Call (RPC) package. (The message protocol is specified with the External Data Representation (XDR) language. See the *External Data Representation Standard: Protocol Specification* Appendix 3 for the details. Here, we assume that the reader is familiar with XDR and do not attempt to justify it or its uses.) The paper by Birrell and Nelson [1] is recommended as an excellent background to and justification of RPC.

Terminology

This chapter discusses servers, services, programs, procedures, clients, and versions. A server is a piece of software where network services are implemented. A network service is a collection of one or more remote programs. A remote program implements one or more remote procedures; the procedures, their parameters, and results are documented in the specific program's protocol specification (see the *Port Mapper Program Protocol* below, for an example). Network clients are pieces of software that initiate remote procedure calls to services. A server may support more than one version of a remote program in order to be forward compatible with changing protocols. For example, a network file service may be composed of two programs. One program may deal with high-level applications such as file system access control and locking. The other may deal with low-level file IO and have procedures like "read" and "write". A client machine of the network file service would call the procedures associated with the two programs of the service on behalf of some user on the client machine.

The RPC Model

The remote procedure call model is similar to the local procedure call model. In the local case, the caller places arguments to a procedure in some well-specified location (such as a result register). It then transfers control to the procedure, and eventually gains back control. At that point, the results of the procedure are extracted from the well-specified location, and the caller continues execution.

The remote procedure call is similar, in that one thread of control logically winds through two processes – one is the caller's process, the other is a server's process. That is, the caller process sends a call message to the server process and waits (blocks) for a reply message. The call message contains the procedure's parameters, among other things. The reply message contains the procedure's results, among other things. Once the reply message is received, the results of the procedure are extracted, and the caller's execution is resumed.

On the server side, a process is dormant awaiting the arrival of a call message. When one arrives, the server process extracts the procedure's parameters, computes the results, sends a reply message, and then awaits the next call message.

Note that in this model, only one of the two processes is active at any given time. However, this model is only given as an example. The RPC protocol makes no restrictions on the concurrency model implemented, and others are possible. For example, an implementation may choose to have RPC calls be asynchronous, so that the client may do useful work while waiting for the reply from the server. Another possibility is to have the server create a task to process an incoming request, so that the server can be free to receive other requests.

Transports and Semantics

The RPC protocol is independent of transport protocols. That is, RPC does not care how a message is passed from one process to another. The protocol deals only with specification and interpretation of messages.

It is important to point out that RPC does not try to implement any kind of reliability and that the application must be aware of the type of transport protocol underneath RPC. If it knows it is running on top of a reliable transport such as TCP/IP[6], then most of the work is already done for it. On the other hand, if it is running on top of an unreliable transport such as UDP/IP[7], it must implement its own retransmission and time-out policy as the RPC layer does not provide this service.

Because of transport independence, the RPC protocol does not attach specific semantics to the remote procedures or their execution. Semantics can be inferred from (but should be explicitly specified by) the underlying transport protocol. For example, consider RPC running on top of an unreliable transport such as UDP/ IP. If an application retransmits RPC messages after short timeouts, the only thing it can infer if it receives no reply is that the procedure was executed zero or more times. If it does receive a reply, then it can infer that the procedure was executed at least once.

A server may wish to remember previously granted requests from a client and not regrant them in order to insure some degree of execute-at-most-once semantics. A server can do this by taking advantage of the transaction ID that is packaged with every RPC request. The main use of this transaction is by the client RPC layer in matching replies to requests. However, a client application may choose to reuse its previous transaction ID when retransmitting a request. The server application, knowing this fact, may choose to remember this ID after granting a request and not regrant requests with the same ID, thus achieving some degree of execute-at-mostonce semantics. The server is not allowed to examine this ID in any other way except as a test for equality.

On the other hand, if using a reliable transport such as TCP/ IP, the application can infer from a reply message that the procedure was executed exactly once, but if it receives no reply message, it cannot assume the remote procedure was not executed. Note that even if a connection-oriented protocol like TCP is used, an application still needs time-outs and reconnection to handle server crashes.

There are other possibilities for transports besides datagramor connection-oriented protocols. For example, a request-reply protocol such as VMTP[2] is perhaps the most natural transport for RPC. At Sun, RPC is currently implemented on top of both TCP/ IP and UDP/IP transports.

Binding and Rendezvous Independence

The act of binding a client to a service is NOT part of the remote procedure call specification. This important and necessary function is left up to some higher-level software. (The software may use RPC itself—see the *Port Mapper Program Protocol* below).

Implementors should think of the RPC protocol as the jumpsubroutine instruction ("JSR") of a network; the loader (binder) makes JSR useful, and the loader itself uses JSR to accomplish its task. Likewise, the network makes RPC useful, using RPC to accomplish this task.

Authentication

The RPC protocol provides the fields necessary for a client to identify itself to a service and vice-versa. Security and access control mechanisms can be built on top of the message authentication. Several different authentication protocols can be supported. A field in the RPC header indicates which protocol is being used. More information on specific authentication protocols can be found in the Authentication Protocols below.

A3.3 RPC Protocol Requirements

The RPC protocol must provide for the following:

- Unique specification of a procedure to be called.
- Provisions for matching response messages to request messages.
- Provisions for authenticating the caller to service and vice-versa.

In addition to these requirements, features that detect the following are worth supporting because of protocol roll-over errors, implementation bugs, user error, and network administration:

- RPC protocol mismatches.
- Remote program protocol version mismatches.
- Protocol errors (such as misspecification of a procedure's parameters).
- · Reasons why remote authentication failed.
- Any other reasons why the desired procedure was not called.

Programs and Procedures

The RPC call message has three unsigned fields: remote program number, remote program version number, and remote procedure number. The three fields uniquely identify the procedure to be called. Program numbers are administered by some central authority (like Sun). Once an implementor has a program number, he can implement his remote program; the first implementation would most likely have the version number of 1. Because most new protocols evolve into better, stable, and mature protocols, a version field of the call message identifies which version of the protocol the caller is using. Version numbers make speaking old and new protocols through the same server process possible.

The procedure number identifies the procedure to be called. These numbers are documented in the specific program's protocol specification. For example, a file service's protocol specification may state that its procedure number 5 is "read" and procedure number 12 is "write".

Just as remote program protocols may change over several versions, the actual RPC message protocol could also change. Therefore, the call message also has in it the RPC version number, which is always equal to two for the version of RPC described here.

The reply message to a request message has enough information to distinguish the following error conditions:

- The remote implementation of RPC does speak protocol version 2. The lowest and highest supported RPC version numbers are returned.
- The remote program is not available on the remote system.
- The remote program does not support the requested version number. The lowest and highest supported remote program version numbers are returned.
- The requested procedure number does not exist. (This is usually a caller side protocol or programming error.)
- The parameters to the remote procedure appear to be garbage from the server's point of view. (Again, this is usually caused by a disagreement about the protocol between client and service.)

Authentication

Provisions for authentication of caller to service and viceversa are provided as a part of the RPC protocol. The call message has two authentication fields, the credentials and verifier. The reply message has one authentication field, the response verifier. The RPC protocol specification defines all three fields to be the following opaque type:

```
auth_flavor
enum
                        {
                         = 0.
        AUTH_NULL
        AUTH_UNIX
                         = 1,
        AUTH_SHORT = 2.
        AUTH_DES
                         = 3
        /* and more to be defined */
};
struct
        opaque_auth {
        auth_flavor flavor;
        opaque body\langle 400 \rangle;
};
```

In simple English, any *opaque_auth* structure is an *auth_flavor* enumeration followed by bytes which are opaque to the RPC protocol implementation.

The interpretation and semantics of the data contained within the authentication fields is specified by individual, independent authentication protocol specifications. (See Authentication Protocols below, for definitions of the various authentication protocols.)

If authentication parameters were rejected, the response message contains information stating why they were rejected.

Program Number Assignment

Program numbers are given out in groups of 0x20000000 (decimal 536870912) according to the following chart:

| Program Number | Description |
|-------------------------|-----------------------------|
| 0x00000000 - 0x1FFFFFFF | Defined by Sun Microsystems |
| 0x20000000 - 0x3FFFFFFF | Defined by User |
| 0x40000000 - 0x5FFFFFFF | Transient |
| 0x60000000 - 0x7FFFFFFF | Reserved |
| 0x80000000 - 0x9FFFFFFF | Reserved |
| 0xA0000000 - 0xBFFFFFFF | Reserved |
| 0xC0000000 - 0xDFFFFFFF | Reserved |
| 0xE0000000 - 0xFFFFFFFF | Reserved |

The first group is a range of numbers administered by Sun Microsystems and should be identical for all sites. The second range is for applications peculiar to a particular site. This range is intended primarily for debugging new programs. When a site develops an application that might be of general interest, that application should be given an assigned number in the first range. The third group is for applications that generate program numbers dynamically. The final groups are reserved for future use, and should not be used.

Other Uses of the RPC Protocol

The intended use of this protocol is for calling remote procedures. That is, each call message is matched with a response message. However, the protocol itself is a message-passing protocol with which other (non-RPC) protocols can be implemented. Sun currently uses, or perhaps abuses, the RPC message protocol for the following two (non-RPC) protocols: batching (or pipelining) and broadcast RPC. These two protocols are discussed but not defined below.

Batching. Batching allows a client to send an arbitrarily large sequence of call messages to a server; batching typically uses reliable byte stream protocols (like TCP/IP) for its transport. In the case of batching, the client never waits for a reply from the server, and the server does not send replies to batch requests. A sequence of batch calls is usually terminated by a legitimate RPC in order to flush the pipeline (with positive acknowledgement).

Broadcast RPC. In broadcast RPC-based protocols, the client sends a broadcast packet to the network and waits for numerous replies. Broadcast RPC uses unreliable, packet-based protocols (like UDP/ IP) as its transports. Servers that support broadcast protocols only respond when the request is successfully processed, and are silent in the face of errors. Broadcast RPC uses the Port Mapper RPC service to achieve its semantics. See the *Port Mapper Program Protocol* below, for more information.

A3.4 The RPC Message Protocol

This section defines the RPC message protocol in the XDR data description language. The message is defined in a top-down style.

```
{
enum
       msg_type
                     = 0.
        CALL
        REPLY
                     = 1
};
* A reply to a call message can take on two forms:
* The message was either accepted or rejected.
*/
enum reply_stat {
  MSG\_ACCEPTED = 0,
  MSG_DENIED
                     = 1
};
/*
* Given that a call message was accepted, the following is the
* status of an attempt to call a remote procedure.
*/
enum accept_stat {
                            = 0, /* RPC executed successfully */
        SUCCESS
                            = 1, /* remote hasn't exported program */
        PROG_UNAVAIL
        PROG_MISMATCH = 2, /* remote can't support version # */
                            = 3, /* program can't support procedure */
        PROC_UNAVAIL
                            = 4 /* procedure can't decode params */
        GARBAGE_ARGS
};
```

```
/*
 * Reasons why a call message was rejected:
 */
enum reject_stat {
        RPC_MISMATCH = 0, /* RPC version number ! = 2 * /
                           = 1 /* remote can't authenticate caller */
        AUTH_ERROR
};
/*
 * Why authentication failed:
 */
enum auth_stat {
                                   = 1, /* bad credentials */
        AUTH_BADCRED
        AUTH_REJECTEDCRED = 2, /* client must begin new session */
        AUTH_BADVERF
                                   = 3, /* bad verifier */
        AUTH_REJECTEDVERF = 4, /* verifier expired or replayed */
                                   = 5 / * rejected for security reasons */
        AUTH_TOOWEAK
};
/*
* The RPC message:
* All messages start with a transaction identifier, xid,
 * followed by a two-armed discriminated union. The union's
 * discriminant is a msg_type which switches to one of the two
 * types of the message. The xid of a REPLY message always
 * matches that of the initiating CALL message. NB: The xid
 * field is only used for clients matching reply messages with
 * call messages or for servers detecting retransmissions; the
 * service side cannot treat this xid as any type of sequence
 * number.
 */ .
struct rpc_msg {
        unsigned int xid;
        union switch (msg_type mtype) {
                case CALL:
                     call_body cbody;
                     case REPLY:
                        reply_body rbody;
        } body;
};
/*
* Body of an RPC request call:
* In version 2 of the RPC protocol specification, rpcvers must
* be equal to 2. The fields prog, vers, and proc specify the
* remote program, its version number, and the procedure within
* the remote program to be called. After these fields are two
```

* authentication parameters: cred (authentication credentials)

```
* and verf (authentication verifier). The two authentication
 * parameters are followed by the parameters to the remote
 * procedure, which are specified by the specific program
* protocol.
 */
struct call_body {
                                 /* must be equal to two (2) */
        unsigned int rpcvers;
        unsigned int prog;
        unsigned int vers;
        unsigned int proc:
        opaque_auth cred;
        opaque_auth verf;
        /* procedure specific parameters start here */
};
/*
* Body of a reply to an RPC request:
* The call message was either accepted or rejected.
*/
union reply_body switch (reply_stat stat) {
        case MSG_ACCEPTED:
                accepted_reply areply;
        case MSG_DENIED:
                rejected_reply rreply;
} reply;
/*
* Reply to an RPC request that was accepted by the server:
 * there could be an error even though the request was accepted.
* The first field is an authentication verifier that the server
 * generates in order to validate itself to the caller. It is
 * followed by a union whose discriminant is an enum
 * accept_stat. The SUCCESS arm of the union is protocol
* specific. The PROG_UNAVAIL, PROC_UNAVAIL, and GARBAGE_ARGP
* arms of the union are void. The PROG_MISMATCH arm specifies
 * the lowest and highest version numbers of the remote program
* supported by the server.
*/
struct accepted_reply {
        opaque_auth verf:
        union switch (accept_stat stat) {
                case
                         SUCCESS:
                         opaque results[0];
                         /* procedure-specific results start here */
                         PROG_MISMATCH:
                case
                         struct {
```

```
unsigned int low:
                                unsigned int high;
                        } mismatch_info;
                default:
                         * Void. Cases include PROG_UNAVAIL.
                        PROC_UNAVAIL,
                         * and GARBAGE_ARGS.
                         */
                        void:
        } reply_data;
};
/*
* Reply to an RPC request that was rejected by the server:
* The request can be rejected for two reasons: either the
* server is not running a compatible version of the RPC
* protocol (RPPC_MISMATCH), or the server refuses to
* authenticate the caller (AUTH_ERROR). In case of an RPC
* version mismatch, the server returns the lowest and highest
* supported RPC version numbers. In case of refused
* authentication, failure status is returned.
*/
union rejected_reply switch (reject_stat stat) {
       case RPC_MISMATCH:
                struct {
                        unsigned int low;
                        unsigned int high:
                } mismatch_info:
       case AUTH_ERROR:
                auth_stat stat;
```

};

A3.5 Authentication Protocols

As previously stated, authentication parameters are opaque, but open-ended to the rest of the RPC protocol. This section defines some "flavors" of authentication implemented at (and supported by) Sun. Other sites are free to invent new authentication types, with the same rules of flavor number assignment as there is for program number assignment.

Null Authentication

Often calls must be made where the caller does not know who he is or the server does not care who the caller is. In this case, the flavor value (the discriminant of the *opaque_auth* 's union) of the RPC message's credentials, verifier, and response verifier is *AUTH_NULL*. The bytes of the opaque_auth's body are undefined. It is recommended that the opaque length be zero.

UNIX Authentication

The caller of a remote procedure may wish to identify himself as he is identified on a UNIX system. The value of the credential's discriminant of an RPC call message is *AUTH_UNIX*. The bytes of the credential's opaque body encode the following structure:

The *stamp* is an arbitrary ID which the caller machine may generate. The *machinename* is the name of the caller's machine (like "krypton"). The *uid* is the caller's effective user ID. The *gid* is the caller's effective group ID. The *gids* is a counted array of groups which contain the caller as a member. The verifier accompanying the credentials should be of AUTH_NULL (defined above).

The value of the discriminant of the response verifier received in the reply message from the server may be AUTH_NULL or AUTH_SHORT. In the case of AUTH_SHORT, the bytes of the response verifier's string encode an opaque structure. This new opaque structure may now be passed to the server instead of the original AUTH_UNIX flavor credentials. The server keeps a cache which maps shorthand opaque structures (passed back by way of an AUTH_SHORT style response verifier) to the original credentials of the caller. The caller can save network bandwidth and server cpu cycles by using the new credentials.

The server may flush the shorthand opaque structure at any time. If this happens, the remote procedure call message will be rejected due to an authentication error. The reason for the failure will be AUTH_REJECTEDCRED. At this point, the caller may wish to try the original AUTH_UNIX style of credentials.

DES Authentication

UNIX authentication suffers from two major problems:

- The naming is too UNIX-system oriented.
- There is no verifier, so credentials can easily be faked.

DES authentication attempts to fix these two problems.

Naming. The first problem is handled by addressing the caller by a simple string of characters instead of by an operating system specific integer. This string of characters is known as the "netname" or network name of the caller. The server is not allowed to interpret the contents of the caller's name in any other way except to identify the caller. Thus, netnames should be unique for every caller in the internet.

It is up to each operating system's implementation of DES authentication to generate netnames for its users that insure this uniqueness when they call upon remote servers. Operating systems already know how to distinguish users local to their systems. It is usually a simple matter to extend this mechanism to the network. For example, a UNIX user at Sun with a user ID of 515 might be assigned the following netname: "unix.515@sun.com". This netname contains three items that serve to insure it is unique. Going backwards, there is only one naming domain called "sun.com" in the internet. Within this domain, there is only one UNIX user with user ID 515. However, there may be another user on another operating system, for example VMS, within the same naming domain that, by coincidence, happens to have the same user ID. To insure that these two users can be distinguished we add the operating system name. So one user is "unix.515@sun.com" and the other is "vms.515@sun.com".

The first field is actually a naming method rather than an operating system name. It just happens that today there is almost a one-to-one correspondence between naming methods and operating systems. If the world could agree on a naming standard, the first field could be the name of that standard, instead of an operating system name.

DES Authentication Verifiers. Unlike UNIX authentication, DES authentication does have a verifier so the server can validate the client's credential (and vice-versa). The contents of this verifier are primarily an encrypted timestamp. The server can decrypt this timestamp, and if it is close to what the real time is, then the client must have encrypted it correctly. The only way the client could encrypt it correctly is to know the "conversation key" of the RPC session. And if the client knows the conversation key, then it must be the real client.

The conversation key is a DES [5] key which the client generates and notifies the server of in its first RPC call. The conversation key is encrypted using a public key scheme in this first transaction. The particular public key scheme used in DES authentication is Diffie-Hellman [3] with 192-bit keys. The details of this encryption method are described later.

The client and the server need the same notion of the current time in order for all of this to work. If network time synchronization cannot be guaranteed, then the client can synchronize with the server before beginning the conversation, perhaps by consulting the Internet Time Server (TIME[4]).

The way a server determines if a client timestamp is valid is somewhat complicated. For any other transaction but the first, the server simply checks for two things:

- the timestamp should be greater than the one previously seen from the same client.
- the timestamp should not have expired.

A timestamp is expired if the server's time is later than the sum of the client's timestamp plus what is known as the client's "window". The "window" is a number the client passes (encrypted) to the server in its first transaction. You can think of it as a lifetime for the credential.

This explains everything but the first transaction. In the first transaction, the server checks only that the timestamp has not expired. However, if this was all that was done, then it would be quite easy for the client to send random data in place of the timestamp with a fairly good chance of succeeding. As an added check, the client sends an encrypted item in the first transaction known as the "window verifier" which must be equal to the window minus 1, or the server will reject the credential.

The client too must check the verifier returned from the server to be sure it is legitimate. The server sends back to the client the encrypted timestamp it received from the client, minus one second. If the client gets anything different than this, it will reject it.

Nicknames and Clock Synchronization. After the first transaction, the server's DES authentication subsystem returns in its verifier to the client an integer "nickname" which the client may use in its further transactions instead of passing its netname, encrypted DES key and window every time. The nickname is most likely an index into a table on the server which stores for each client its netname, decrypted DES key and window.

Although they originally were synchronized, the client's and server's clocks can get out of sync again. When this happens the client RPC subsystem most likely will get back *RPC_AUTHER*-*ROR*, at which point it should resynchronize.

A client may still get the *RPC_AUTHERROR* error even though it is synchronized with the server. The reason is that the server's nickname table is a limited size, and it may flush entries whenever it wants. A client should resend its original credential in this case and the server will give it a new nickname. If a server crashes, the entire nickname table gets flushed, and all clients will have to resend their original credentials.

/*

DES Authentication Protocol (in XDR language)

^{*} There are two kinds of credentials: one in which the client uses

^{*} its full network name, and one in which it uses its "nickname"

^{* (}just an unsigned integer) given to it by the server. The

^{*} client must use its fullname in its first transaction with the

^{*} server, in which the server will return to the client its

^{*} nickname. The client may use its nickname in all further

^{*} transactions with the server. There is no requirement to use the

^{*} nickname, but it is wise to use it for performance reasons.

```
enum authdes_namekind {
        ADN_FULLNAME = 0,
        ADN_NICKNAME = 1
};
* A 64-bit block of encrypted DES data
typedef opaque des_block[8];
/*
* Maximum length of a network user's name
const MAXNETNAMELEN = 255;
/*
* A fullname contains the network name of the client, an encrypted
* conversation key and the window. The window is actually a
* lifetime for the credential. If the time indicated in the
* verifier timestamp plus the window has past, then the server
* should expire the request and not grant it. To insure that
* requests are not replayed, the server should insist that
* timestamps are greater than the previous one seen, unless it is
* the first transaction. In the first transaction, the server
* checks instead that the window verifier is one less than the
* window.
*/
struct authdes_fullname {
        string name (MAXNETNAMELEN); /* name of client */
                                           /* PK encrypted conversation key */
        des_block key;
                                           /* encrupted window */
        unsigned int window;
};
* A credential is either a fullname or a nickname
union authdes_cred switch (authdes_namekind adc_namekind) {
        case ADN_FULLNAME:
                authdes_fullname adc_fullname;
        case ADN_NICKNAME:
                unsigned int adc_nickname;
};
/*
* A timestamp encodes the time since midnight, January 1, 1970.
*/
```

```
struct timestamp {
        unsigned int seconds: /* seconds */
        unsigned int useconds: /* and microseconds */
};
/*
* Verifier: client variety
* The window verifier is only used in the first transaction. In
* conjunction with a fullname credential, these items are packed
* into the following structure before being encrypted:
* struct {
* adv_timestamp:
                          -one DES block
* adc_fullname.window; -one half DES block
* adv_winverf:
                          -one half DES block
* }
* This structure is encrypted using CBC mode encryption with an
* input vector of zero. All other encryptions of timestamps use
* ECB mode encryption.
*/
struct authdes_verf_clnt {
        timestamp adv_timestamp; /* encrypted timestamp */
        unsigned int adv_winverf; /* encrypted window verifier */
};
/*
* Verifier: server variety
* The server returns (encrypted) the same timestamp the client
* gave it minus one second. It also tells the client its nickname
* to be used in future transactions (unencrypted).
*/
struct authdes_verf_svr {
        timestamp adv_timeverf;
                                     /* encrypted verifier */
        unsigned int adv_nickname; /* new nickname for client */
};
```

Diffie-Hellman Encryption. In this scheme, there are two constants, *BASE* and *MODULUS*. The particular values Sun has chosen for these for the DES authentication protocol are:

```
const BASE = 3;
const MODULUS = "d4a0ba0250b6fd2ec626e7efd637df76c716e22d0944b88b"; /*
hex */
```

The way this scheme works is best explained by an example. Suppose there are two people "A" and "B" who want to send encrypted messages to each other. So, A and B both generate "secret" keys at random which they do not reveal to anyone. Let these keys be represented as SK(A) and SK(B). They also publish in a public directory their "public" keys. These keys are computed as follows:

PK(A) = (BASE ** SK(A)) mod MODULUS PK(B) = (BASE ** SK(B)) mod MODULUS

The "**" notation is used here to represent exponentiation. Now, both A and B can arrive at the "common" key between them, represented here as CK(A, B), without revealing their secret keys.

A computes:

 $CK(A, B) = (PK(B) ** SK(A)) \mod MODULUS$

while B computes:

 $CK(A, B) = (PK(A) ** SK(B)) \mod MODULUS$

These two can be shown to be equivalent:

(PK(B) ** SK(A)) mod MODULUS = (PK(A) ** SK(B)) mod MODULUS

We drop the "mod MODULUS" parts and assume modulo arithmetic to simplify things:

PK(B) ** SK(A) = PK(A) ** SK(B)

Then, replace PK(B) by what B computed earlier and likewise for PK(A).

((BASE ** SK(B)) ** SK(A) = (BASE ** SK(A)) ** SK(B)

which leads to:

BASE ** (SK(A) * SK(B)) = BASE ** (SK(A) * SK(B))

This common key CK(A, B) is not used to encrypt the timestamps used in the protocol. Rather, it is used only to encrypt a conversation key which is then used to encrypt the timestamps. The reason for doing this is to use the common key as little as possible, for fear that it could be broken. Breaking the conversation key is a far less serious offense, since conversations are relatively short-lived.

The conversation key is encrypted using 56-bit DES keys, yet the common key is 192 bits. To reduce the number of bits, 56 bits are selected from the common key as follows. The middle-most 8bytes are selected from the common key, and then parity is added to the lower order bit of each byte, producing a 56-bit key with 8 bits of parity.

A3.6 Record Marking Standard

When RPC messages are passed on top of a byte stream protocol (like TCP/IP), it is necessary, or at least desirable, to delimit one message from another in order to detect and possibly recover from user protocol errors. This is called record marking (RM). Sun uses this RM/TCP/IP transport for passing RPC messages on TCP streams. One RPC message fits into one RM record.

A record is composed of one or more record fragments. A record fragment is a four-byte header followed by 0 to $(2^{**}31) - 1$ bytes of fragment data. The bytes encode an unsigned binary number; as with XDR integers, the byte order is from highest to lowest. The number encodes two values—a boolean which indicates whether the fragment is the last fragment of the record (bit value 1 implies the fragment is the last fragment), and a 31-bit unsigned binary value which is the length in bytes of the fragment's data. The boolean value is the highest-order bit of the header; the length is the 31 low-order bits. (Note that this record specification is NOT in XDR standard form!)

A3.7 The RPC Language

Just as there was a need to describe the XDR data-types in a formal language, there is also need to describe the procedures that operate on these XDR data-types in a formal language as well. We use the RPC Language for this purpose. It is an extension to the XDR language. The following example is used to describe the essence of the language.

An Example Service Described in the RPC Language

Here is an example of the specification of a simple ping program.

```
/*
* Simple ping program
*/
program PING_PROG {
/* Latest and greatest version */
version PING_VERS_PINGBACK {
void
PINGPROC_NULL(void) = 0;
```

```
* Ping the caller, return the round-trip time
        * (in microseconds). Returns -1 if the operation
        * timed out.
        */
        int
        PINGPROC_PINGBACK(void) = 1;
\} = 2:
/*
* Original version
*/
version PING_VERS_ORIG {
        void
        PINGPROC_NULL(void) = 0;
        \} = 1;
\} = 1;
const PING_VERS = 2:
                               /* latest version */
```

The first version described is *PING_VERS_PINGBACK* with two procedures, *PINGPROC_NULL* and *PINGPROC_PINGBACK*. *PINGPROC_NULL* takes no arguments and returns no results, but is useful for computing round-trip times from the client to the server and back again. By convention, procedure 0 of any RPC protocol should have the same semantics, and never require any kind of authentication. The second procedure is used for the client to have the server do a reverse ping operation back to the client, and it returns the amount of time (in microseconds) that the operation used. The next version, *PING_VERS_ORIG*, is the original version of the protocol and it does not contain *PINGPROC_ PINGBACK* procedure. It is useful for compatibility with old client programs, and as this program matures it may be dropped from the protocol entirely.

The RPC Language Specification

The RPC language is identical to the XDR language, except for the added definition of a *program-def* described below.

```
program-def:

"program" identifier "{"

version-def

version-def *

"}" "=" constant ";"
```

```
version-def:

"version" identifier "{"

procedure-def

procedure-def *

"}" "=" constant ";"
```

```
procedure-def:
```

```
type-specifier identifier "(" type-specifier ")" "=" constant ";"
```

Syntax Notes

- The following keywords are added and cannot be used as identifiers: "program" and "version";
- A version name cannot occur more than once within the scope of a program definition. Nor can a version number occur more than once within the scope of a program definition.
- A procedure name cannot occur more than once within the scope of a version definition. Nor can a procedure number occur more than once within the scope of version definition.
- Program identifiers are in the same name space as constant and type identifiers.
- Only unsigned constants can be assigned to programs, versions and procedures.

A3.8 Port Mapper Program Protocol

The port mapper program maps RPC program and version numbers to transport-specific port numbers. This program makes dynamic binding of remote programs possible.

This is desirable because the range of reserved port numbers is very small and the number of potential remote programs is very large. By running only the port mapper on a reserved port, the port numbers of other remote programs can be ascertained by querying the port mapper.

The port mapper also aids in broadcast RPC. A given RPC program will usually have different port number bindings on different machines, so there is no way to broadcast directly to all of these programs. The port mapper, however, does have a fixed port

number. So, to broadcast to a given program, the client actually sends its message to the port mapper located at the broadcast address. Each port mapper that picks up the broadcast then calls the local service specified by the client. When the port mapper gets the reply from the local service, it sends the reply on back to the client.

```
Port Mapper Protocol Specification
(in RPC Language)
```

```
const PMAP_PORT = 111;
                                   /* portmapper port number */
/*
* A mapping of (program, version, protocol) to port number
*/
struct mapping {
        unsigned int prog;
        unsigned int vers;
        unsigned int prot;
        unsigned int port;
};
/*
* Supported values for the "prot" field
*/
                                /* protocol number for TCP/IP */
const IPPROTO_TCP = 6;
const IPPROTO_UDP = 17;
                                    /* protocol number for UDP/IP */
/*
* A list of mappings
*/
struct *pmaplist {
        mapping map;
        pmaplist next;
};
/*
* Arguments to callit
*/
struct call_args {
        unsigned int prog;
        unsigned int vers;
        unsigned int proc;
        opaque args\langle\rangle;
};
```

```
* Results of callit
*/
struct call_result {
        unsigned int port;
        opaque res\langle \rangle;
};
/*
* Port mapper procedures
*/
program PMAP_PROG {
        version PMAP_VERS {
                void
                PMAPPROC_NULL(void)
                                                 = 0;
                bool
                PMAPPROC_SET(mapping)
                                                 = 1:
                bool
                PMAPPROC_UNSET(mapping)
                                                 = 2;
                unsigned int
                PMAPPROC\_GETPORT(mapping) = 3;
                pmaplist
                PMAPPROC_DUMP(void)
                                                 = 4;
               call_result
                PMAPPROC_CALLIT(call_args)
                                                 = 5;
        \} = 2;
\} = 100000;
```

Port Mapper Operation

The portmapper program currently supports two protocols (UDP/IP and TCP/IP). The portmapper is contacted by talking to it on assigned port number 111 (SUNRPC [8]) on either of these protocols. The following is a description of each of the portmapper procedures:

PMAPPROC_NULL:

This procedure does no work. By convention, procedure zero of any protocol takes no parameters and returns no results.

PMAPPROC_SET:

When a program first becomes available on a machine, it registers itself with the port mapper program on the same machine. The program passes its program number "prog", version number "vers", transport protocol number "prot", and the port "port" on which it awaits service request. The procedure returns a boolean response whose value is *TRUE* if the procedure successfully established the mapping and *FALSE* otherwise. The procedure refuses to establish a mapping if one already exists for the tuple "(prog, vers, prot)".

PMAPPROC_UNSET:

When a program becomes unavailable, it should unregister itself with the port mapper program on the same machine. The parameters and results have meanings identical to those of *PMAPPROC*_*SET*. The protocol and port number fields of the argument are ignored.

PMAPPROC_GETPORT:

Given a program number "prog", version number "vers", and transport protocol number "prot", this procedure returns the port number on which the program is awaiting call requests. A port value of zero means the program has not been registered. The "port" field of the argument is ignored.

PMAPPROC_DUMP:

This procedure enumerates all entries in the port mapper's database. The procedure takes no parameters and returns a list of program, version, protocol, and port values.

PMAPPROC_CALLIT:

This procedure allows a caller to call another remote procedure on the same machine without knowing the remote procedure's port number. It is intended for supporting broadcasts to arbitrary remote programs via the well-known port mapper's port. The parameters "prog", "vers", "proc", and the bytes of "args" are the program number, version number, procedure number, and parameters of the remote procedure. Note:

- 1. This procedure only sends a response if the procedure was successfully executed and is silent (no response) otherwise.
- 2. The port mapper communicates with the remote program using UDP/IP only.

The procedure returns the remote program's port number, and the bytes of results are the results of the remote procedure.

References

[1] Birrell, Andrew D. & Nelson, Bruce Jay; "Implementing Remote Procedure Calls"; XEROX CSL-83-7, October 1983.

[2] Cheriton, D.; "VMTP: Versatile Message Transaction Protocol", Preliminary Version 0.3; Stanford University, January 1987.

[3] Diffie & Hellman; "Net Directions in Cryptography"; IEEE Transactions on Information Theory IT-22, November 1976.

[4] Harrenstien, K.; "Time Server", RFC 738; Information Sciences Institute, October 1977.

[5] National Bureau of Standards; "Data Encryption Standard"; Federal Information Processing Standards Publication 46, January 1977.

[6] Postel, J.; "Transmission Control Protocol—DARPA Internet Program Protocol Specification", RFC 793; Information Sciences Institute, September 1981.

[7] Postel, J.; "User Datagram Protocol", RFC 768; Information Sciences Institute, August 1980.

[8] Reynolds, J. & Postel, J.; "Assigned Numbers", RFC 923; Information Sciences Institute, October 1984.

Differences Between the RPC Library on SunOS 4.0 and SunOS 4.1 Appendix

4

This appendix covers the differences between the SunOS 4.0 and the SunOS 4.1 versions of the RPC Library, the *portmap* service, and the *rpcgen* program. Only a brief summary of the changes are presented. More detailed information can be obtained from the SunOS 4.1 documentation.

A4.1 RPC Library

The following new features were added to the RPC Library:

- The routine *xdrrec_readbytes()* was added to the Library. This routine can only be used on streams created by *xdrrec_create()*. The routine attempts to read a specified number of bytes from the XDR stream into a specified buffer.
- The routine *clntudp_bufcreate()* is now documented and available to the user. This routine is in the SunOS 4.0 version of the Library but is for internal use by the Library. This routine is the same as *clntudp_create()* except that you can now specify the size of the send and receive buffers.
- The routine *clnt_create_vers()* was added to the Library. This routine is a generic client creation routine which also checks for the version available. Remember that the *clnt_create()* routine returns a valid client handle even if the specified version number supplied to the routine is not registered with the *portmap* service. However, *clnt_create_vers()* does this for you and returns a valid handle only if a version within the range supplied is supported by the server.

• New request values have been added to the *clnt_control()* routine. The new request values, their associated argument types, and what the requests do follows:

| CLGET_FD | int | get socket descriptor associated with the client handle |
|-----------------|------|---|
| CLSET_FD_CLOSE | void | close socket when <i>clnt_destroy()</i> |
| CLSET_FD_NCLOSE | void | leave socket open when <i>clnt_destroy()</i> is called |

• You need to include $\langle rpc/raw.h \rangle$ when using the raw transport creation routines.

The following fixes were added to the SunOS 4.1 version of the Library:

• The routine *clnt_call()* now does an exponential back off when retrying RPC requests.

- The *clnt_broadcast()* routine now uses an exponential back off on the retry time-out. In SunOS 4.0 version, it used a linear series for retry time-outs.
- The *svc_run()* routine now ignores all error codes returned from *select()* except for EBADF. The *svc_run()* in the SunOS 4.0 version of the Library ignores only the EINTR error code. This change was made because EBADF is the only error that should result in *svc_run()* returning. This change also makes *svc_run()* more tolerant of applications that handle signals and inadvertently modify the global variable *errno* inside their signal handling routine. The problem exists in *svc_run()* because the return from *select()* and the test of the *errno* variable is not an atomic operation and, in fact, control may be passed to a signal handling routine during this sequence of operations. The SunOS 4.1 version of the *svc_run()* routine was covered in Chapter 6.
- The raw transport creation routines have been fixed.

A4.2 Portmap Service

The following fixes were made to the *portmap* service:

• Disallows PMAP_SET and PMAP_UNSET operations from remote hosts.

- Disallows PMAP_SET and PMAP_UNSET of reserved ports from non-reserved ports.
- The PMAP_RMTCALL procedure no longer forks. This provides a significant performance improvement for the processing of broadcast requests.

A4.3 Rpcgen

The following features were added to *rpcgen*:

- Can generate servers which can be invoked by the *inetd* program.
- Allows for -DDEFINE statements on the command line to define macros.
- Allows for server error messages to be logged using the syslog mechanism.
- Can generate an indexed-by-procedure table.

Source for Examples

The examples used in this book can be obtained in a computer readable form. To order the examples, specify the type of media desired and the format to be used on the disk. The media, media formats supported, and prices are listed below. Appendix 5

| 3.5″ | Floppy Disk UNIX tar Format | \$8.00(US) |
|-------|--------------------------------|-------------|
| 5.25″ | Floppy Disk UNIX tar Format | \$8.00(US) |
| 1/4″ | Streaming Tape Unix tar Format | \$12.00(US) |
| 3.5″ | Floppy Disk MS-DOS Format | \$8.00(US) |
| 5.25″ | Floppy Disk MS-DOS Format | \$8.00(US) |
| | | |

Send your order along with payment to:

RPC Programming Examples P.O. Box 12474 El Paso, TX 79912

Bibliography

Bergan, E.S., Tolchin, S.G.: Using Remote Procedure Calls (RPC) for a Distributed Clinical Information System. Proceedings of the UniForum Conference, 1986

Birman, K.P., Joseph, T.A.: Exploiting Replication in Distributed Systems. In: Mullender, S.(ed.): Distributed Systems. Reading (Mass.): Addison-Wesley 1989, pp. 191 - 214

Birrell, A.D., Nelson, B.J.: Implementing Remote Procedure Calls. ACM Transactions on Computer Systems 2:1, 39 - 59 (1984)

Comer, D.: Internetworking With TCP/IP -- Principles, Protocols and Architectures. Englewood Cliffs (New Jersey): Prentice Hall 1988

Coulouris, G.F., Dollimore, J.: Distributed Systems -- Concepts and Design. Reading (Mass.): Addison-Wesley 1988

Diffie, W., Hellman, M.E.: New Directions in Cryptography. IEEE Transactions on Information Theory *IT*-22, 644 - 654 (1976)

Folts, H.: Open Systems Standards -- OSI Remote and Reliable Operations. IEEE Network 3:3, (1989)

Havender, J.W.: Avoiding Deadlock in Multitasking Systems. IBM Systems Journal 7:2, 74 - 84 (1968)

Leffler, S.J., et al.: The Design and Implementation of the 4.3BSD UNIX Operating System. Reading (Mass.): Addison-Wesley 1989

Lyon, B.: Sun Remote Procedure Call Specification. Sun Microsystems, Inc., Mountain View (Calif.), 1984.

Lyon, B.: Sun External Data Representation Specification. Sun Microsystems, Inc., Mountain View (Calif.), 1984.

National Bureau of Standards: Data Encryption Standard. Federal Information Processing Standards Publication 46, January 15, 1977

Needham, R.M., Schroeder, M.D.: Using Encryption for Authentication in Large Networks of Computers. Communications of the ACM 21:12, (1978)

Obermarck, R.: Distributed Deadlock Detection Algorithm. ACM Transactions on Database Systems 7:2, 187 - 208 (1982)

Samar, V., McManis, C.: Sun Remote Procedure Call Implementation Made Transport Independent. Sun Microsystems, Mountain View (Calif.), December 1989.

Spector, A.Z.: Distributed Transaction Processing Facilities. In: Mullender, S.(ed.): Distributed Systems. Reading (Mass.): Addison-Wesley 1989, pp. 191 - 214

Stallings, W.: Handbook of Computer-Communications Standards (Volume 1). New York (New York): Macmillan 1987

Sun Microsystems: Network Programming. Mountain View (Calif.): Sun Microsystems 1988

Sun Microsystems: SunOS Reference Manual. Mountain View (Calif.): Sun Microsystems 1988

Sun Microsystems: Security Features Guide. Mountain View (Calif.): Sun Microsystems 1988

Taylor, B., Goldberg, D.: Secure Networking in the Sun Environment. Proceedings of the USENIX Summer Conference, 1986

Weihl, W.E.: Remote Procedure Call. In: Mullender, S.(ed.): Distributed Systems. Reading (Mass.): Addison-Wesley 1989, pp. 65 - 86

White, J.E.: A High-level Framework For Network-based Resource Sharing. Proceedings of the National Computer Conference, June 1976

Xerox Corporation: Courier: The Remote Procedure Call Protocol -- Xerox System Integration Standard XSIS-038112. Stamford (Conn.): Xerox Corporation 1981

Index

Abstract Syntax Notation 1 (ASN.1), 258 API (application programming interface), 9 Application programming interface, see **API** entries Applications, distributed, see Distributed applications ASN.1 (Abstract Syntax Notation 1), 250 Asynchronous RPC, 147-158 Authentication fields, 286-287 Authentication mechanism, 115, 116-126 Authentication protocols, 56–60, 291– 298Authentication types, 56–57, 115–126 DES, 57-60, 121-126, 293-298 null, 57, 116–117, 292 UNIX, 57, 117–120, 292–293 Batch-mode RPC, 162–170, 287–288 Big-endian byte order, 12–14 Binding independence, 284 Breakpoints, setting, 243 Broadcast RPC, 147, 158-162, 288 Broken connections, 220-221 Byte order, 12-14 Bytes, 14

Cache replies, 217–222 CAD (computer-aided design), 5 Call messages, 51, 53–54, 282, 288–291 Callback deadlocks, 222–223 Callback RPC, 147, 149–158 *Callit* procedure, 63–64 CK (conversation key), 58–60 Client crash, 7, 211 Client handle, 93, 102–106, 221, 225, 235 Client-server model, 6–7 Client side low-level RPC routines, 102 - 106Client stub interface, 179 Client timestamp, 294-295 clnt stat values, 112, 113-115 Code-generation tool, 255 Command line options to rpcgen, 193-194 Computer-aided design (CAD), 5 Computer names, 65 Computer networks, see Network entries Computer readable sources for examples, 311 Connections, broken, 220-221 Conversation key (CK), 58-60, 298 "Cooked credentials," 118 Courier, 10 Credentials, 56-57, 115, 286 "cooked," 118 null, 57, 116 "raw," 118 Credential window values, 131 Data base routines (DBM), 126-144 Data-communication networks, see Network entries Data description, 11 Data Encryption Standard, see DES entries Data representation, 11 eXternal, see XDR entries XDR, 2, 4, 11-45, 66, 261-280 Datagram transport, 89, 92–93, 212–215 DBM (data base routines), 126–144 de facto and de jure standards, 247-248 Deadlocks, callback, 222-223

DES (Data Encryption Standard), 57 DES authentication, 120-126, 293-298 DES authentication protocol, 295-297 DES authentication verifiers, 294-295 Development tools for distributed applications, 254-256 DFS (distributed file systems), 247-248 Diffie-Hellman encryption, 297-298 Dispatch procedure, 100 server. see Server dispatch routine Distributed applications, 4-6 callback deadlocks in, 222-223 design issues in, 212-229 designing, 207-212 developing RPC-based, 207-243 development tools for, 254-256 error recovery in, 211, 216-222 execute-at-most-once semantics in, 215 - 216idempotency in, 215-216 implementing, 8 make file, 211 partitioning, 210 performance in, 223-229 purpose of, 208 testing procedures for, 229-243 testing with network component, 236, 242 - 243testing without network component, 230 - 236time-outs in, 229 transport issues in, 212-225 Distributed file systems (DFS), 247-248 Distributed systems, 1-2 8086 architecture, 12–14 Encryption, Diffie-Hellmen, 297–298 Error recovery in distributed applications, 211, 216-222 Execute-at-most-once semantics, 215-216Explicit typing, 250 eXternal Data Representation routines, see XDR entries Fault tolerance, 5 File handle, 224 File state, 7

Grade-reporter programs, 71-88 rpcgen version, 181–193 revised, 93-101 Group ID, 121 Headers, 22-23 Heterogeneity issue, 209 Host names, 65 Hyper integers, 264 Idempotency, 215-216 Identification field, transaction (XID), 53 - 54IEEE (Institute of Electrical and Electronics Engineers), 248 Institute of Electrical and Electronics Engineers (IEEE), 248 Interface, 68 application programming (API), 9 client stub, 179 programming, 4 transport, 91-93 International Organization for Standardization (ISO), 248-250 Internet Protocol (IP), 92 IP (Internet Protocol), 92 ISO (International Organization for Standardization), 248-250 keuserv daemon, 121

Library dbm, 126–144 RPC, see RPC Library Transport Interface, 92 XDR, see XDR Library entries Linking client and server, 231–235 Little-endian byte order, 12–14 Local data representations, 15 Lyon, Bob, 10

Magic number, 41 Manycast RPC, 246 Maximum transfer unit (MTU), 158 MTU (maximum transfer unit), 158 Multicast feature, 229

Name server, 60 Netnames, 120-125, 293 NETPATH variable, 252 Network access. 1-2 Network clients, 281; see also Client entries Network component testing with, 236, 242-243 testing without, 230-236 Network File System (NFS), 10 Network resources, 1 Network service, 9, 66-67, 281 Network-transport dependency, 89, 211 NFS (Network File System), 10 Nicknames, 295 Nonblocking RPC, 147-149 Null authentication, 57, 116, 292 Null credential, 57, 117 Null verifiers, 117 Octets, 14 ONC (Open Network Computing) platform, 2 Opaque data, 32-33 **Open Network Computing (ONC) plat**form. 2 **Open Systems Interconnection (OSI)** model, 248-249 Operating-system dependency, 211 **OSI** (Open Systems Interconnection) Application Layer, 249 Data Link Layer, 249 model, 248-249 Network Layer, 249 Physical Layer, 249 Presentation Layer, 249, 250 Session Layer, 249 Transport Layer, 249 P1003.8 working group, 250-251

Partitioning distributed applications, 210
Performance analysis tool, 255–256
Performance issue, 210, 223–229
ping program, 241
Platform, 2
Pointer declarations, 204–205

Portable Operating Systems Environment Committee (POSIX), 248, 250 - 251portmap program, 61, 62, 242 procedures, 62-64 portmap service, fixes made to, 308-309 Portmapper, 61-62 Port map protocol, 301–305 Portmapper Library Routines, 175–177 pmap getmaps(), 176 pmap_getport(), 176-177 pmap_rmtcall(), 177 pmap_set(), 176-177 pmap_unset(), 177 Ports, 60-61 **POSIX** (Portable Operating Systems Environment Committee), 248, 250-251Preprocessor directives for rpcgen, 194-195Procedure, 65 remote, see Remote Procedure entries Protocol specification, 179, 210-211 Protocol specification file, 198 Protocol testers, 243 Prototyping tool, 254–255 "Raw credentials," 119 Raw transport, 235-236, 237-241 References, 280, 305 Release independence, 209-210 Reliability, 52 Remote computer, 3 **Remote Operations Service Element** (ROSE), 250 Remote procedure, 3, 65 Remote Procedure Call, see RPC entries Remote procedure identification, 66-69 Remote program, 281 Remote program definitions, 197–199 Remote program numbers, 68 Reply cache, 217-222 Reply messages, 51, 54–58 RFC1014, 261 RFC1057, 281 **ROSE** (Remote Operations Service Element), 250

RPC (Remote Procedure Call), 2 API standard, 247 argument encoding and decoding, 66 asynchronous, 147-158 batch-mode, 162-170 broadcast, 147, 158-162, 288 callback, 147, 149-158 history, 9-10 ISO and, 248-250 manycast, 246 message authentication, 284; see also Authentication entries message-passing scheme, 3, 51-53 nonblocking, 147-149 POSIX and, 250-251 procedure numbers, 68-69, 285 program number assignment, 287 program versions, 67 multiple, supporting, 81, 88 **RPC** Administrator, 258 RPC-based distributed applications, see **Distributed** applications RPC history, 9–10 RPC information, Sun, 257-259 RPC Language, see RPCL entries RPC Library, 3-4, 69-70 additional features, 147-177 differences between SunOS 4.0 and SunOs 4.1, 307-309 future directions of, 245-247 new features added to, 307-308 processing errors in, 109, 110–112 RPC mechanism, 65 RPC message protocol, 288-291 RPC model, 282 RPC program numbers, 257–259 RPC programming, 65-89 future directions of, 245-256 high-level, 70–88 low-level, 91-145 **RPC** Programming Examples, 311 RPC protocol requirements, 284-288 RPC protocol specification, 281–305 **RPC** Library routines high-level, 79-81 callrpc(), 80-81registerrpc(), 79-80svc_run(), 80

low-level client side, 102-106 auth_destroy(), 116 authnone_create(), 116 authdes_create(), 121 authunix_create(), 117 authunix_create_default(), 118 clnt_broadcast(), 159-160 $clnt_call(), 105$ clnt control(), 106 clnt_create(), 102-103 clnt_destroy(), 105 clnt_freeres(), 106 clnt_pcreateerror(), 112 clnt_perrno(), 112 clnt_perror(), 113 clnt_spcreateerror(), 112 clnt_sperrno(), 113 clnt sperror(), 113 clntraw_create(), 235 clnttcp_create(), 104-105 clntudp create(), 103-104 example, 126-145 server side, 98-102 svc_destroy(), 99 svc_freeargs(), 101-102 svc_getargs(), 101 svc_getreqset(), 173 svc_register(), 100-101 $svc_sendreply(), 102$ svc_unregister(), 101-102 svcerr_auth(), 110 svcerr_decode(), 110 svcerr_noproc(), 110 svcerr_noprog(), 111 svcerr_progvers(), 111 svcerr_systemerr(), 111 svcerr_weakauth(), 111 svcraw_create(), 99 svctcp_create(), 99 svcudp_create(), 98-99 support, 122-123 host2netname(), 122 getnetname(), 123 netname2host(), 122-123 netname2user(), 123 user2netname(), 123 RPC source code availability, 257

RPC standardization efforts, 247-251 RPC version number, 54 rpcbind program, 253-254 rpcgen program, 8, 179-206 command line options to, 193-194 features added to, 309 preprocessor directives for, 194-195 rpcinfo utility, 230 RPCL (RPC Language), 179, 195-206, 299 - 301arrays, fixed links, 269 Character, 226 fixed length, 202, 268-269 string, 237 variable length, 202-203, 269 Boolean type, 202, 264 comment delimiter, 196 constants, 201, 277, 278 data types, 199-201, 262 discriminated union, 34, 203-204, 270-271.278enumerations, 201-202, 263-264, 277 fixed-length arrays, 202, 268-270, 277, 285fixed-length opaque data, 266–267 floating-point data, 264–276 double-precision, 265-266, 278 single-precision, 265, 278 identifiers, 196, 276 integers, 262, 263 hyper, 264 unsigned, 263 unsigned hyper, 264 opaque data, 32-33, 213 fixed-length, 205, 266-267 variable-length, 205, 267-268 procedure definition syntax, 197–198 program definition syntax, 197-198 reserved keywords, 196, 278 strings, 31-32, 205, 227, 268 structures, 203, 270, 277 typedefs, 205-206, 271-272, 278 types, decomposition of, 199-200 typing, explicit, 250 union, discriminated, 34 unsigned hyper integers, 264 unsigned integers, 263 variable-length arrays, 202-203, 269

variable-length opaque data, 205, 267–268 variable names, 278 version definition syntax, 197 void declaration, 201, 271 RPCL specification, 300–301 RPCL syntax notes, 301 RPCL unions, 203–204 RPCSRC source code, 4, 257 *rwhod* daemon, 228

Security issues, 208-209 Server, 6, 281 linking client and, 231–235 stateless and stateful, 6-7 Server crash, 7, 218-222 Server dispatch routine, 107–109 example, 109, 110-111 Server process, 6–7 Server side low-level RPC routines, 99-103 Server skeleton, 179 Service request structure svc rea(), 116 Session Layer, 249 Simulation tool, 255 Sockets, 91-92 SPARC architecture, 12–14 Standard data representations, 15 Standardization efforts, RPC, 243-251 Stateless and stateful servers, 6–7 STREAMS input/output mechanism. 251Sun Microsystems, 2 Sun Operating System, see SunOS entries Sun RPC information, 257-259 SunOS (Sun Operating System), 4 SunOS 4.0 versus SunOS 4.1, 307-309 svc_run() routine, writing, 170–175 SVR4 (System V Release 4), 248 Syntax information, XDR, 277-279 Syntax RPCL, 277-278, 300-301 XDR, 278-279 System/370 architecture, 12–14 System V Release 4 (SVR4), 248

TCP (Transmission Control Protocol), 92, 212-214, 283-284 Testers, protocol, 243 Testing procedures for distributed applications, 229-243 Testing tool, 255 TI RPC (transport-independent RPC), 248, 251-254 Time-outs in distributed applications. 229Timestamp, client, 294–295 TLI (Transport Layer Interface), 251 Transaction identification (XID) field, 53 - 54Transfer unit, maximum (MTU), 158 Transmission Control Protocol (TCP), 98, 212-214, 283-284 Transport, 4 raw, 235-236, 237-241 Transport address, 60 Transport dependencies, 209 Transport handle, 93 Transport-independent RPC (TI RPC), 248, 251-254 Transport interface, 91-93 **Transport Interface Library**, 92 Transport Layer Interface (TLI), 251 Transport protocols, 283–284 Transport selection mechanism, 251-253**TRUE**. 66 Type-def syntax information, 278 UDP (User Datagram Protocol), 89, 92, 212-214, 283-284 UNIX. 2 UNIX authentication, 117-120, 292-293 UNIX-style credential, 57 User Datagram Protocol (UDP), 89, 92, 212 - 214**User ID**, 120

VAX architecture, 12–14 Vendor dependencies, 210 Verifiers, 56–57, 116 DES authentication, 294–295

null, 116 **UNIX. 117** Version definition syntax, 197 White, Jim, 10 Wrapper, 78-79 XDR (eXternal Data Representation), 2, 4, 11-49, 66, 261-280 data types, 199-206, 262-274 in-line macros, 29, 227 macros, 29 in-line, 29, 227 stream-regulated, 26-27 memory allocation, 39-41 memory freeing, 40-41 memory streams, 20-22 example, 21 record stream, 22-26 example, 41, 45, 46-49 record marking (RM), 299 record streams, 22-26 standard I/O streams, 19-20 example, 41-45stream-related macros, 26, 27 XDR block size, 16–17, 225–227, 262, 274, 275 XDR data description example, 279–280 XDR data representation, 15–17 XDR definition, 14–15 XDR filters, 17-18, 27-41 composite, 31-38 custom, 38-39 primitive, 28-31 XDR handles, 18 XDR language specification, 276–279 XDR lexical notes, 276-277 **XDR** Library routines $xdr_array(), 33-34$ xdr_bytes(), 32-33 $xdr_char(), 28$ xdr_destroy(), 27 $xdr_double(), 28$ $xdr_enum(), 28$ $xdr_float(), 28$ $xdr_free(), 40$

 $xdr_getpos(), 26$ xdr_inline(), 29 xdr_int(), 28 $xdr_long(), 28$ xdrmem_create(), 20-22 xdr_opaque(), 32 $xdr_pointer(), 36-37$ xdrrec_create(), 22-25 xdrrec_endofrecord(), 25 xdrrec_eof(), 26 xdrrec_skiprecord(), 26 xdr_reference(), 35-36 xdr_setpos(), 27 $xdr_short(), 28$ xdrstdio_create(), 19-20 $xdr_string(), 31-32$ *xdr_u_char()*, 28 $xdr_u(int(), 28$ $xdr_u_long(), 28$

 $xdr_u_short(), 28$ xdr_union(), 34-35 $xdr_vector(), 33$ $xdr_void(), 28$ xdr_wrapstring(), 32 XDR standard, 11, 261-280 areas for future enhancement, 273-274 data types in, 262-274 defined, 261-262 status of, 261 XDR streams, 18-27 XDR syntax information, 277-278 XDR syntax notes, 278-279 XDR unit size, 274, 275 XID (transaction identification) field, 53 - 54

ypbind process, 228-229