

12. Naming

Any problem in computing can be solved by another level of indirection.
David Wheeler

Introduction

This handout is about orderly ways of naming complicated collections of objects in a computer system. A basic technique for understanding a big system is to describe it as a collection of simple parts. Being able to name these parts is a necessary aspect of such a description, and often the most important aspect.

The basic idea can be expressed in two ways that are more or less equivalent:

Identify values by variable length names called *path names* that are sequences of simple names that are strings. Think of all the names with the same prefix (for instance, `/udir/lampson` and `/udir/lynch`) as being grouped together. This grouping induces a tree structure on the names. Non-leaf nodes in the tree are *directories*.

Make a tree of nodes with simple names on the arcs. The leaf nodes are values and the internal nodes are *directories*. A node is named by a path through the tree from the root; such a name is called a *path name*.

Thus `/udir/lampson/pocs/handouts/12` is a path name for a value (perhaps the text of this handout), and `/udir/lampson/pocs/handouts` is a path name for a directory (other words for directory are folder, context, closure, environment, binding, and dictionary). The collection of all the path names that make sense in some situation is called a *name space*. Viewing a name space as a tree gives us the standard terminology of parents, children, ancestors, and descendants.

Using path names to name values (or objects, if you prefer) is often called ‘hierarchical naming’ or ‘tree-structured naming’. There are a lot of other names for it that are used in special situations: mounting, search paths, multiplexing, device addressing, network references. An important reason for studying naming in general is that you don’t have to start from scratch in understanding all those other things.

Path names are good because:

- The name space can grow indefinitely, and the growth can be managed in a decentralized way. That is, the authority to create names in one part of the space can be delegated, and thereafter there is no need for synchronization. Names that start `/udir/lampson` are independent of names that start `/udir/rinard`.
- Many kinds of data can be encapsulated under this interface, with a common set of operations. Arbitrary operations can be encoded as reads and writes of suitably chosen names.

As we have seen, a path name is a sequence of simple names. We use the types $N = \text{string}$ for a simple name and $PN = \text{SEQ } N$ for a path name. It is often convenient to write a path name as a string. The syntax of these strings is not important; it is just a convention for encoding the path names. Here are some examples:

```
/udir/lampson/pocs/handouts/12
lampson@mediaone.net
```

```
16.23.5.193
```

Unix path name

Internet mail address. The path name is
{ "net", "mediaone", "lampson" }

IP network address (fixed length)

We will normally write path names as Unix file names, rather than as the sequence constructors that would be correct Spec. Thus `a/b/c/1026` instead of $PN\{ "a", "b", "c", "1026" \}$.

People often try to distinguish a name (what something is) from an address (where it is) or a route (how to find it). This is a matter of levels of abstraction and must not be taken as absolute. At a given level of abstraction we tend to identify objects at that level by names, the lower-level objects that code them by addresses, and paths at lower levels by routes. Examples:

```
microsoft.com -> 207.46.130.149 -> SEQ [router output port, LAN address]
a/b/c/1026 -> INode/1026 -> DA/2 -> [cylinder, head, sector, byte 2]
```

Sometimes people talk about “descriptive names”, which are queries in a database. We will see that these are readily encompassed within the framework of path names. That is a formal relationship, however. There is an important practical difference between a *designator* for a single entity, such as `lampson@mediaone.net`, and a *description* or query such as “everyone at MIT’s LCS whose research involves parallel computing”. The difference is illuminated by the comparison between the name `eecsfaculty@eecs.mit.edu` and the query “the faculty members in MIT’s EECS department”. The former name is probably maintained with some care; it’s anyone’s guess how reliable the answer to the query is. When using a name, it is wise to consider whether it is a designator or a description.

This is not to say that descriptions or queries are bad. On the contrary, they are very valuable, as any one knows who has ever used a web search engine. However, they usually work well only when a person examines the results with some care.

In the remainder of this handout we examine the specs for the two ways of describing a name space that we introduced earlier: as a memory addressed by path names, and as a tree (or more generally a graph) of directories. The two ways are closely related, but they give rise to somewhat different specs. Then we study the recursive structure of name spaces and various ways of inducing a name space on a collection of values. This leads to a more abstract analysis of how the spec for a name space can vary, depending on the properties of the underlying values. We conclude our general treatment by examining how to name a name space. Finally, we give a large number of examples of name spaces; you might want to look at these first to get some more context.

Name space as memory

We can view a name space as an example of the memory abstraction we studied earlier. Recall that a memory is a partial map $M = A \rightarrow V$. Here we take $A = PN$ and replace M with D (for

directory). This kind of memory differs from the byte-addressable physical memory of a computer in several ways¹:

- The map is partial.
- The domain is changing.
- The current value of the domain (that is, which names are defined) is interesting.
- `PN`'s with the same prefix are related (though not as much as in the second view of name spaces).

Here are some examples of name spaces that can naturally be viewed as memories:

The Simple Network Management Protocol (SNMP) is used to manage components of the Internet. It uses path names (rooted in IP addresses) to name values, and the basic operations are to read and write a single named value.

Several file systems use a single large table to map the path name of a file to the extents that represent it.

```
MODULE MemNames0 [V] EXPORT Read, Write, Remove, Enum, Next, Rename =
```

```
TYPE N          = String          % Name
   PN           = SEQ N WITH {"<<=":=PNLE} % Path Name
   D            = PN -> V         % Directory
VAR d           := D{}            % the state
FUNC PNLE(pn1, pn2) -> Bool = pn1.LexLE(pn2, N."<=") % pn1 <= pn2
```

Here are the familiar `Read` and `Write` procedures; `Read` raises `error` if `d` is undefined at `pn`, for consistency with later specs. In this basic spec none of the other procedures raises `error`; this innocence will not persist when things get more complicated. It's common to also have a `Remove` procedure for making a `PN` undefined; note that unlike a file system, this `Remove` does not erase the values of longer names that start with `PN`. This is because, unlike a file system, this spec does not ensure that every prefix of a defined `PN` is defined.

```
FUNC Read(pn) -> V RAISES {error} = RET d(pn) [*] RAISE error
APROC Write(pn, v) = << d := d{pn -> v} >>
APROC Remove(pn) = << d := d{pn -> } >>
```

The body of `Write` is usually written `d(pn) := v`.

It's important that the map is partial, and that the domain changes. This means that we need operations to find out what the domain is. Simply returning the entire domain is not practical, since it may be too big, and usually only part of it is of interest. There are two schools of thought

¹ It differs much less from the virtual memory, in which the map may be partial and the domain may change as new virtual memory is assigned or files are mapped. Actually these things can happen to physical memory as well, especially in the part of it implemented by I/O devices.

about what form these operations should take, represented by the functions `Enum` and `Next`; only one of these is needed.

`Enum` returns all the simple names that can lead to a value starting from `pn`; another way of saying this is that it returns all the names bound in the directory named `pn`. By recursively applying `Enum` to `pn + n` for each simple name `n` that `Enum` returns, you can explore the entire tree.

On the other hand, if you keep feeding `Next` its own output, starting with `{}`, it walks the tree of defined names depth-first, returning in turn each `PN` that is bound to a `v`. It finishes with `{}`.

Note that what `Next` does is not the same as returning the results of `Enum` one at a time, since `Next` explores the entire tree, not just one directory. Thus `Enum` takes the organization of the name space into directories more seriously than does `Next`.

```
FUNC Enum(pn) -> SET N = RET {pn1 | d!(pn + pn1) | pn1.head}
FUNC Next(pn) -> PN = VAR later := {pn' | d!pn' /\ pn <= pn'} |
RET later.fmin(PN."<=") [*] RET {} % {} if later is empty
```

A separate issue is arranging to get a reasonable number of results from one of these procedures. If the directory is large, `Enum` as defined here may return an inconveniently large set, and we may have to call `Next` inconveniently many times. In real life we would make either routine return a sequence of `N`'s or `PN`'s, usually called a 'buffer'. This is a standard use of batching to reduce the overhead of invoking an operation, without allowing the batches to get too large. We won't add this complication to our specs.

Finally, there is a `Rename` procedure that takes directories quite seriously. It reflects the idea that all the names which start the same way are related, by changing all the names that start with `from` so that they start with `to`. Because directories are not very real in the representation, this procedure has to do a lot of work. It erases everything that starts with either argument, and then copies everything in the original `d` that starts with `from` to the corresponding path name that starts with `to`. `Read x <= y` as "x is a prefix of y".

```
APROC Rename(from: PN, to: PN) RAISES {error} = << VAR d0 := d |
IF from <= to => RAISE error % can't rename to a descendant
[*] DO VAR pn :IN d.dom | (to <= pn \/ from <= pn) => d := d{pn -> } OD;
DO VAR pn | d(to + pn) # d0(from + pn) => d(to + pn) := d0(from + pn) OD
FI >>
END MemNames0
```

Here is a different version of `Rename` that makes explicit the relation between the initial state `d` and the final state `d'`. `Read x >= y` as "x is a suffix of y".

```
APROC Rename(from: PN, to: PN) RAISES {error} = <<
IF VAR d' |
  (ALL x: PN, y: PN | (
    x >= from => ~ d'!x
    [*] x = to + y /\ d!(from + y) => d'(x) = d(from + y)
    [*] ~ x >= to /\ d!x => d'(x) = d(x)
    [*] ~ d'!x )
  )
```

```

=> d := d'
[*] RAISE error FI >>

```

There is often a rule that a name can be bound to a directory or to a value, but not both. For this we need a slightly different spec that marks a name as bound to a directory by giving it the special value `isD`, with a separate procedure for making an empty directory. To enforce the new rule every routine can now raise `error`, and `Remove` erases the whole sub-tree. As usual, `boxes` mark the changes from `MemNames0`.

```

MODULE MemNames [V] EXPORT Read, Write, MakeD, Remove, Enum, Rename =

TYPE Dir = ENUM[isDir]
D = PN -> (V + Dir) [SUCHTHAT (\d | d({}) IS Dir)] % root a Dir
VAR d := D[{} -> isDir]

% INVARIANT (ALL pn, pn' | d!pn' /\ pn' > pn => d(pn) = isDir)
FUNC Read(pn) -> V RAISES {error} = [d(pn) IS V =>] RET d(pn) [*] RAISE error
FUNC Enum(pn) -> SET N RAISES {error} =
  [d(pn) IS Dir =>] RET {pn1 | d!(pn + pn1) | pn1.head} [*] RAISE error
APROC Write(pn, v) RAISES {error} = << [Set(pn, v)] >>
APROC MakeDir(pn) RAISES {error} = << [Set(pn, isDir)] >>

APROC Remove(pn) = % Erase everything with pn prefix.
  << [DO VAR pn' :IN d.dom | (pn <= pn') =>] d := d{pn' ->} [DD] >>

APROC Rename(from: PN, to: PN) RAISES {error} = << VAR d0 := d |
  IF from <= to => RAISE error % can't rename to a descendant
  [*] DO VAR pn :IN d.dom | (to <= pn /\ from <= pn) => d := d{pn ->} OD;
  DO VAR pn | d(to + pn) # d0(from + pn) =>
    d(to + pn) := d0(from + pn) OD
  FI >>

APROC Set(pn, y: (V + D) RAISES {error} =
  << pn # {} /\ d(pn.reml) IS D =>] d(pn) := y [*] RAISE error >>

END MemNames

```

A file system usually forbids overwriting a file with a directory (for no obvious reason) or overwriting a non-empty directory with anything (because a directory is precious and should not be clobbered wantonly), but these rules are rather arbitrary, and we omit them here.

Exercise: write a version of `Rename` that makes explicit the relation between the initial state `d` and the final state `d'`, in the style of the second `Rename` of `MemNames0`.

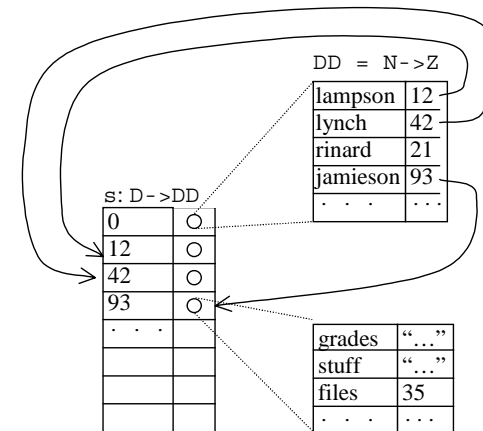
The `MemNames` spec is basically the same as the simple `Memory` spec. Complications arise because the domain can change, and because of the distinction between directories and values. The specs in the next section take this distinction much more seriously.

Name space as graph of directory objects

These specs are reasonably simple, but they are clumsy for operations on directories such as `Rename`. More fundamentally, they don't handle aliasing, where the same object has more than one name. The other (and more usual) way to look at a hierarchical name space is to think of each directory as a function that maps a simple name (not a path name) to a value or another directory, rather than thinking of the entire tree as a single `PN -> V` map. This tree (or general graph) structure maps a `PN` by mapping each `N` in turn, traversing a path through the graph of directories; hence the term 'path name'. We continue to use the type `D` for a directory.

Our eventual goal is a spec for a name space as graph that is 'object-oriented' in the sense that you can supply different code for each directory in the name space. We will begin, however, with a simpler spec that is equivalent to `MemNames`, evolve this to a more general spec that allows aliases, and finally add the object orientation.

The obvious thing to do is to make a `D` be a function `N -> Z`, where `Z = (D + V)` as before, and have a state variable `d` which is the root of the tree. Unfortunately this completely functional structure doesn't work smoothly, because there's no way to change the value of `a/b/c/d` without changing the value of `a/b/c` so that it contains the new value of `a/b/c/d`, and similarly for `a/b` and `a` as well.²



We solve this problem in the usual way with another level of indirection, so that the value of a directory name is not a `N -> Z` but some kind of reference or pointer to a `N -> Z`, as shown in the figure. This reference is an 'internal name' for a directory. We use the name `DD` for the actual

² The method of explicitly changing all the functions up to the root has some advantages. In particular, we can make several changes to different parts of the name space appear atomically by waiting to rewrite the root until all the changes are made. It is not very practical for a file system, though at least one has been built this way: H.E. Sturgis, *A Post-Mortem for a Time-sharing System*, PhD thesis, University of California, Berkeley, and Report CSL 74-1, Xerox Research Center, Palo Alto, Jan 1974. It has also been used in database systems to atomically change the entire database state; in this context it is called 'shadowing'. See Gray and Reuter, pp 728-732.

function $N \rightarrow Z$ and introduce a state variable s that holds all the DD values; its type is $D \rightarrow DD$. A D is just the internal name of a directory, that is, an index into s . We take $D = \text{Int}$ for simplicity, but any type with enough values would do; in Unix $D = \text{INO}$. You may find it helpful to think of D as a pointer and s as a memory, or of D as an inode number and s as the inodes. Later sections explore the meaning of a D in more detail, and in particular the meaning of root .

Once we have introduced this extra indirection the name space does not have to be a tree, since two PN 's can have the same D value and hence refer to the same directory. In a Unix file system, for example, every directory with the path name pn also has the path names $pn/.$, $pn/./.$, etc., and if pn/a is a subdirectory, then the parent also has the names $pn/a/.$, $pn/a/./a/.$, etc. Thus the name space is not a tree or even a DAG, but a graph with cycles, though the cycles are constrained to certain stylized forms involving $'.'$ and $'..'$. This means, of course, that there are defined PN 's of unbounded length; in real life there is usually an arbitrary upper bound on the length of a defined PN .

The spec below does not expose D 's to the client, but deals entirely in PN 's. Real systems often do expose the D pointers, usually as some kind of capability (for instance in a file system that allows you to open a directory and obtain a file descriptor for it), but sometimes just as a naked pointer (for instance in many distributed name servers). The spec uses an internal function Get , defined near the end, that looks up a PN in a directory; GetD is a variation that raises error if it can't return a D .

MODULE ObjNames0 [V] EXPORT Read, Write, MakeD, Remove, Enum, Rename =

```

TYPE D          = Int                % just an internal name
    Z          = (V + D)            % the value of a name
    DD         = N -> Z             % a Directory

CONST root      : D := 0
VAR s           := (D -> DD){}{root -> DD{}} % initially empty root

FUNC Read(pn) -> V RAISES {error} = VAR z := Get(root, pn) |
    IF z IS V => RET z [*] RAISE error FI

FUNC Enum(pn) -> SET PN RAISES {error} = RET s(GetD(root, pn)).dom
% Raises error if pn isn't a directory, like MemNames.

A write operation on the name a/b/c has to change the d component of the directory a/b; it does
this through the procedure SetPN, which gets its hands on that directory by invoking
GetD(root, pn.reml).

APROC Write(pn, v) RAISES {error} = << SetPN(pn, v) >>
APROC MakeD(pn) RAISES {error} = << VAR d := NewD() | SetPN(pn, d) >>

APROC Remove(pn) RAISES {error} =
    << VAR d := GetD(root, pn.reml) | >>

APROC Rename(from: PN, to: PN) RAISES {error} = <<
    IF (to = {}) \ / (from <= to) => RAISE error % can't rename to a descendant
    [*] VAR fd := GetD(root, from.reml), % know from, to # {}
        td := GetD(root, to .reml) |
        s(fd)!(from.last) =>

```

```

    s(td) := s(td)(to .last -> s(fd)(from.last));
    s(fd) := s(fd){from.last -> }
[*] RAISE error
FI >>

```

The remaining routines are internal. The main one is $\text{Get}(d, pn)$, which returns the result of starting at d and following the path pn . GetD raises error if it doesn't get a directory. NewD creates a new, empty directory.

```

FUNC Get(d, pn) -> Z RAISES {error} =
% Return the value of pn looked up starting at z.
    IF pn = {} => RET d
    [*] VAR z := s(d)(pn.head) | z IS D => RET Get(z, pn.tail)
    [*] RAISE error
    FI

FUNC GetD(d, pn) -> D RAISES {error} = VAR z := Get(d, pn) |
    IF z IS D => RET z [*] RAISE error FI

APROC SetPN(pn, z) RAISES {error} =
    << VAR d := GetD(root, pn.reml) | s(d)(pn.last) := z >>

APROC NewD() -> D = << VAR d | ~ s!d => s(d) := DD{}; RET d >>

END ObjNames0

```

As we did with the second version of MemNames0 . Rename , we can give a definition of Get in terms of a predicate. It says that there's a sequence p of directories starting at d and ending at the result of Get , such that the components of pn select the corresponding components of p ; if there's no such sequence, raise error .

```

FUNC Child(z1, z2) -> Bool = z1 IS D /\ s!z1 /\ z2 IN s(z1).rng

FUNC Get(d, pn) -> Z RAISES {error} = <<
    IF VAR p :IN Child.paths |
        p.head = d /\ (ALL i :IN pn.dom | p(i+1) = s(p(i)(pn(i)))) => RET p.last
    [*] RAISE error
    FI >>

```

ObjNames0 is equivalent to MemNames . The abstraction function from ObjNames0 to MemNames is

$$\text{MemNames.d} = (\backslash pn \mid G(pn) \text{ IS } V \Rightarrow G(pn) \text{ [*] } G(pn) \text{ IS } D \Rightarrow \text{isD})$$

where we define a function G which is like Get on root except that it is undefined where Get raises error :

```

FUNC G(pn) -> Z = RET Get(root, pn) EXCEPT error => IF false => SKIP FI

```

The **EXCEPT** turns the error exception from Get into an undefined result for G .

Exercise: What is the abstraction function from MemNames to ObjNames0 .

Objects, aliases, and atomicity

This spec makes clear the basic idea of interpreting a path name as a path through a graph of directories, but it is unrealistic in several ways:

The operations for changing the value of the DD functions in `s` may be very different from the `Write` and `Maked` operations of `ObjNames0`. This happens when we impose the naming abstraction on a data structure that changes according to its own rules. `SNMP` is a good example; the values of names changes because of the operation of the network. Later in this handout we will explore a number of these variations.

There is often an ‘alias’ or ‘symbolic link’ mechanism which allows the value of a name `n` in context `d` to be a *link* (`d', pn`). The meaning is that `d(n)` is a synonym for `Get(d', pn)`.

The operations are specified as atomic, but this is often too strong.

Our next spec, `ObjNames`, reflects all these considerations. It is rather complicated, but the complexity is the result of the many demands placed on it; ideas for simplifying it would be gratefully received. `ObjNames` is a fairly realistic spec for a naming system that allows for both symbolic links and extensible code for directories.

A `ObjNames.D` has `get` and `set` methods to allow for different code, though for now we don’t take any advantage of this, but use the fixed code `GetFromS` and `SetInS`. In the section on object-oriented directories below, we will see how to plug in other versions of `D` with different `get` and `set` methods. The section on coherence below explains why `get` is a procedure rather than a function. These methods map undefined values to `nil` because it’s tricky to program with undefined in this general setting; this means that `Z` needs `Null` as an extra case.

`Link` is another case of `Z` (the internal value of a name), and there is code in `Get` to follow links; the rules for doing this are somewhat arbitrary, but follow the Unix conventions. Because of the complications introduced by links, we usually use `GetDN` instead of `Get` to follow paths; this procedure converts a `PN` relative to `root` into a directory `d` and a name `n` in that directory. Then the external procedures read or write the value of that name.

Because `Get` is no longer atomic, it’s no longer possible to define it in terms of a path through the directories that exists at a single instant. The section on atomicity below discusses this point in more detail.

MODULE `ObjNames [V]` EXPORT ... =

```

TYPE D          = Int % Just an internal name
               WITH {get:=GetFromS, set:=SetInS} % get returns nil if undefined
Link           = [d: (D + Null), pn] % d=nil for 'relative': the containing D
Z             = (V + D + Link + Null) % nil means undefined
DD            = N -> Z

CONST root    : D := 0
VAR s         := (D -> DD){}{root -> DD{}} % initially empty root

APROC GetFromS(d, n) -> Z = % d.get(n)
  << RET s(d)(n) [*] RET nil >>

APROC SetInS (d, n, z) = % d.set(n, z)
% If z = nil, SetInS leaves n undefined in s(d).
  << IF z # nil => s(d)(n) := z [*] s(d) := s(d){n -> } FI >>

```

```

PROC Read (pn) -> V RAISES {error} = VAR z := Get(root, pn) |
  IF z IS V => RET z [*] RAISE error FI

PROC Enum (pn) -> SET N RAISES {error} =
% Can't just write RET GetD(root, pn).get.dom as in ObjNames0, because get isn't a function.
% The lack of atomicity is on purpose.
VAR d := GetD(root, pn), ns: SET N := {}, z |
  DO VAR n | << z := d.get(n); ~ n IN ns /\ z # nil => ns + := {n} >> OD;
  RET ns

PROC Write (pn, v) RAISES {error} = SetPN(pn, v, true)

PROC Maked(pn) RAISES {error} = VAR d := NewD() | SetPN(pn, d, false)

PROC Rename(from: PN, to: PN) RAISES {error} = VAR d, n, d', n' |
  IF (to = {}) \ / (from <= to) => RAISE error % can't rename to a descendant
  [*] (d, n) := GetDN(from, false); (d', n') := GetDN(to, false);
  << d.get!n => d'.set(n', d.get(n)); d.set(n, nil) >>
  [*] RAISE error
FI

```

This version of `Rename` imposes a different restriction on renaming to a descendant than real file systems, which usually have a notion of a distinguished parent for each directory and disallow `ParentPN(d) <= ParentPN(d')`. They also usually require `d` and `d'` to be in the same ‘file system’, a notion which we don’t have. Note that `Rename` does its two writes atomically, like many real file systems.

The remaining routines are internal. `Get` follows every link it sees; a link can appear at any point, not just at the end of the path. `GetDN` would be just

```

IF pn = {} => RAISE error [*] RET (GetD(root, pn.reml), pn.last) FI

```

except for the question of what to do when the value of this `(d, n)` is a link. The `followLastLink` parameter says whether to follow such a link or not. Because this can happen more than once, the body of `GetDN` needs to be a loop.

```

PROC Get(d, pn) -> Z RAISES {error} = VAR z := d |
% Return the value of pn looked up starting at d.
  DO << pn # {} => VAR n := pn.head, z' |
    IF z IS D => % must have a value for n.
      z' := z.get(n);
      IF z' # nil =>
        % If there's a link, follow it. Otherwise just look up n.
        IF (z, pn') := FollowLink(z, n); pn := pn' + pn.tail
          [*] z := z' ; pn := pn.tail
        FI
      [*] RAISE error
    FI
  [*] RAISE error
  FI
  >> OD; RET z

PROC GetD(d, pn) -> D RAISES {error} = VAR z := Get(d, pn) |
  IF z IS D => RET z AS D [*] RAISE error FI

```

```

PROC GetDN(pn, followLastLink: Bool) -> (D, N) RAISES {error} = VAR d := root |
% Convert pn into (d, n) such that d.get(n) is the item that pn refers to.
DO IF pn = {} => RAISE error
  [*] VAR n := pn.last, z |
    d := Get(d, pn.reml);
    % If there's a link, follow it and loop. Otherwise return.
    << followLastLink => (d, pn) := FollowLink(d, n) [*] RET (d, n) >>
  FI
OD

```

```

APROC FollowLink(d, n) -> (D, PN) = <<
% Fail if d.get(n) not Link. Use d as the context if the link lacks one.
VAR l := d.get(n) | l IS Link => RET ((l.d IS D => l.d [*] d), l.pn) >>

```

```

PROC SetPN(pn, z, followLastLink: Bool) RAISES {error} =
  VAR d, n | (d, n) := GetDN(pn, followLastLink); d.set(n, z)

APROC NewD() -> D = << VAR d | ~ s!d => s(d) := D{}; RET d >>

END ObjNames

```

Object-oriented directories

Although `D` in `ObjNames` has `get` and `set` methods, they are the same for all `D`'s. To encompass the full range of applications of path names, we need to make a `D` into a full-fledged 'object', in which different instances can have different `get` and `set` operations (yet another level of indirection). This is the essential meaning of 'object-oriented': the type of an object is a record of routine types which defines a single interface to all objects of that type, but every object has its own values for the routines, and hence its own code.

To do this, we change the type to:

```

TYPE D          = [get: APROC (n) -> Z, set: PROC (n, z) RAISES {error}]
DR              = Int                                % what D used to be; R for reference
keeping the other types from ObjNames unchanged:
Z               = (V + D + Link + Null)             % nil means undefined
DD              = N -> Z

```

We also need to change the state to:

```

CONST root      := NewD()
VAR s           := (DR -> DD) {root -> DD{}}        % initially empty root

```

and to provide a new version of the `NewD` procedure for creating a new standard directory. The routines that `NewD` assigns to `get` and `set` have the same bodies as the `GetFromS` and `SetInS` routines.

A technical point: The reason for not writing `get := s(dr)` in `NewD` below is that this would capture the value of `s(dr)` at the time `NewD` is invoked; we want the value at the time `get` is invoked, and this is what we get because of the fact that `Spec` functions are functions on the global state, rather than pure functions.

```

APROC NewD() -> D = << VAR dr | ~ s!dr =>
  s(dr) := DD{};
  RET D{ get := (\ n | s(dr)(n)),

```

```

  set := (PROC (n, z) = IF z # nil => s(dr)(n) := z
    [*] s(dr) := s(dr){n -> } FI) }

```

```

PROC SetErr(n, z) RAISES {error} = RAISE error
% For later use as a set proc if the directory is read-only

```

We don't need to change anything else in `ObjNames`.

We will see many other examples of `get` and `set` routines. Note that it's easy to define a `D` that disallows updates, by making `set` be `SetErr`.

Views and recursive structure

In this section we examine ways of constructing name spaces, and in particular ways of building up directories out of existing directories. We already have a basic recursive scheme that makes a set of existing directories the children of a parent directory. The generalization of this idea is to define a function on some state that returns a `D`, that is, a pair of `get` and `set` procedures. There are various terms for this:

- 'encapsulating' the state,
- 'embedding' the state in a name space,
- 'making the state compatible' with a name space interface,
- defining a 'view' on the state.

We will usually call it a view. The spec for a view defines how the result of `get` depends on the state and how `set` affects the state.

All of these terms express the same idea: make the state behave like a `D`, that is, abstract it as a pair of `get` and `set` procedures. Once packaged in this way, it can be used wherever a `D` can be used. In particular, it can be an argument to one of the recursive views that make a `D` out of other `D`'s: a parent directory, a link, or the others discussed below. It can also be the argument of tools like the Unix commands that list, search, and manipulate directories.

The read operations are much the same for all views, but updates vary a great deal. The two simplest cases are the one we have already seen, where you can set the value of a name just as you write into a memory location, and the even simpler one that disallows updates entirely; the latter is only interesting if `get` looks at global state that can change in other ways, as it does in the `Union` and `Filter` operations below. Each time we introduce a view, we will discuss the spec for updating it.

In the rest of this section we describe views that are based on directories: links, mounting, unions, and filters. The final section of the handout gives many examples of views based on other kinds of data.

Links and mounting

The idea behind links (called 'symbolic links' in Unix, 'shortcuts' in Windows, and 'aliases' in the Macintosh) is that of an alias (another level of indirection): we can define the value of a name in a directory by saying that it is the same as the value of some other name in some other

directory. If the value is a directory, another way of saying this is that we can represent a directory d by the link (d', pn') , with $d(pn) = d'(pn')(pn)$, or more graphically $d/pn = d'/pn'/pn$. When put in this form it is usually called *mounting* the directory $d'(pn')$ on $pn0$, if $pn0$ is the name of d . In this language, $pn0$ is called a ‘mount point’. Another name for it is ‘junction’.

We have already seen code in `ObjNames` to handle links. You might wonder why this code was needed. Why isn’t our wonderful object-oriented interface enough? The reason is that people expect more from aliases than this interface can deliver: there can be an alias for a value, not only for a directory, and there are complicated rules for when the alias should be followed silently and when it should be an object in its own right that can be enumerated or changed

Links and mounting make it possible to give objects the names you want them to have, rather than the ones they got because of defects in the system or other people’s bad taste. A very down-to-earth example is the problems caused by the restriction in standard Unix that a file system must fit on a single disk. This means that in an installation with 4 disks and 12 users, the name space contains `/disk1/john` and `/disk2/mary` rather than the `/udir/john` and `/udir/mary` that we want. By making `/udir/john` be a link to `/disk1/john`, and similarly for the other users, we can hide this annoyance.

Since a link is not just a D , we need extra interface procedures to read the value of a link (without following it automatically, as `Read` does), and to install a link. We call the install procedure `Mount` to emphasize that a mount point and a symbolic link are essentially the same thing. The `Mount` procedure is just like `Write` except for the second argument’s type and the fact that it doesn’t follow a final link in pn .

```
PROC ReadLink(pn) -> Link RAISES {error} = VAR d, n |
  (d, n) := GetDN(pn, false);
  VAR z | z := d.get(n); IF z IS Link => RET z [*] RAISE error FI

PROC Mount(pn, link) -> DD = SetPN(pn, link, false)
```

The section on roots below discusses where we might get the D in the `link` argument of `Mount`. In the common case of a link to someplace in the same name space, we have:

```
PROC MakeLink(pn, pn', local: Bool) =
  Mount(pn, Link{d := (local => nil [*] root), pn := pn'})
```

Updating (with `Write`, for instance) makes sense when there are links, but there are two possibilities. If every link is followed then a link never gets updated, since `GetDN` never returns a reference to a link. If a final link is not followed then it can be replaced by something else.

What is the relation between these links and what Unix calls ‘hard links’? A Unix hard link is an inode number, which you can think of as a direct pointer to a file; it corresponds to a D in `ObjNames`. Several directory entries can have the same inode number. Another way to look at this is that the inodes are just another kind of name of the form `inodeRoot/2387754`, so that a hard link is just a link that happens to be an inode number rather than an ordinary path name. There is no provision for making the value of an inode number be a link (or indeed anything except a file), so that’s the end of the line.

Unions

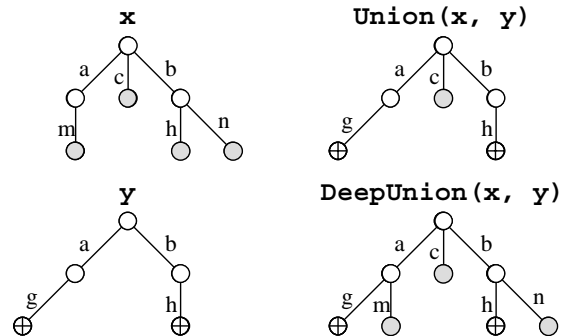
Since a directory is a function $N \rightarrow Z$, it is natural to combine two directories with the “+” overlay operator on functions³. If we do this repeatedly, writing $d1 + d2 + d3$, we get the effect of a ‘search path’ that looks at $d3$ first, then $d2$, and finally $d1$ (in that order because “+” gives preference to its second argument, unlike a search path which gives preference to its first argument). The difference is that this rule is part of the name space, while a search path must be coded separately in each program that cares. It’s unclear whether an update of a union should change the first argument, change the second argument, do something more complicated, or raise an error. We take the last view for simplicity.

```
FUNC Union(d1, d2) -> D = RET D{get := d1.get + d2.get, set := SetErr}
```

Another kind of union combines the name spaces at every level, not just at the top level, by merging directories recursively. This is the most general way to combine two trees that have evolved independently.

```
FUNC DeepUnion(d1, d2) -> D = RET D{
  get := (\ n |
    ( d1.get(n) IS D /\ d2.get(n) IS D => DeepUnion(d1.get(n), d2.get(n))
    [*] (d1.get + d2.get)(n) ),
  set := SetErr}
```

This is a spec, of course, not efficient code.



Filters and queries

Given a directory d , we can make a smaller one by selecting some of d ’s children. We can use any predicate for this purpose, so we get:

```
FUNC Filter(d, p: (D, N) -> Bool) -> D =
  RET D{get := (\ n | (p(d, n) => d.get(n)) [*] nil ), set := SetErr}
```

³ See section 9 of the Spec reference manual.

⁴ This is a bit oversimplified, since `get` is an APROC and hence doesn’t have “+” defined. But the idea should be clear. Plan 9 (see the examples at the end) implements unions.

Examples:

Pattern match in a directory: `a/b/*.ps`. The predicate is true if `n` matches `*.ps`.

Querying a table: `payroll/salary>25000/name`. The predicate is true if `Get(d, n/salary) > 25000`. See the example of viewing a table in the final section of examples.

Full text indexing: `bwl/papers/word:naming`. The predicate is true if `d.get(n)` is a text file that contains the word `naming`. The code could just search all the text files, but a practical one will probably involve an auxiliary index structure that maps words to the files that contain them, and will probably not be perfectly coherent.

See the ‘semantic file system’ example below for more details and a reference.

Variations

It is useful to summarize the ways in which a spec for a name space might vary. The variations almost all have to do with the exact semantics of updates:

What operations are updates, that is, can change the results of `Read`?

Are there *aliases*, so that an update to one object can affect the value of others?

Are the updates *atomic*, or it is possible for reads to see intermediate states? Can an update be lost, or partly lost, if there is a crash?

Viewed as a memory, is the name space *coherent*? That is, does every read that follows an update see the update, or is it possible for the old state to hang around for a while?

How much can the set of defined `DN`'s change? In other words, is it useful to think about a *schema* for the name space that is separate from the current state?

Updates

If the directories are ‘real’, then there will be non-trivial `Write`, `MakeD`, and `Rename` operations. If they are not, these operations will always raise `error`, there will be operations to update the underlying data, and the view function will determine the effects of these updates on `Read` and `Enum`. In many systems, `Read` and `Write` cannot be modeled as operations on memory because `Write(a, r)` does not just change the value returned by `Read(a)`. Instead they must be understood as methods of (or messages sent to) some object.

The earliest example of this kind of system is the DEC Unibus, the prototype for modern I/O systems. Devices on such an I/O bus have ‘device registers’ that are named as locations in memory. You can read and write them with ordinary load and store instructions. Each device, however, is free to interpret these reads and writes as it sees fit. For example, a disk controller may have a set of registers into which you can write a command which is interpreted as “read `n` disk blocks starting at address `da` into memory starting at address `a`”. This might take three writes, for the parameters `n`, `da`, and `a`, and the third write has the side effect of starting execution of the command.

The most recent well-known incarnation of this idea is the World Wide Web, in which read and write actions (called `Get` and `Post` in the protocol) are treated as messages to servers that can search databases, accept orders for pizza, or whatever.

Aliases

We have already discussed this topic at some length. Links and unions both introduce aliases. There can also be ‘hard links’, which are several occurrences of the same `D`. In a Unix file system, for example, it is possible to have several directory entries that point to the same file. A hard link differs from a soft link because the connection it establishes between a name and a file cannot be broken by changing the binding of some other name. And of course a view can introduce arbitrarily complicated aliasing. For example, it’s fairly common for an I/O device that has internal memory to make that memory addressable with two control registers `a` and `v`, and the rule that a read or write of `v` refers to the internal memory location addressed by the current contents of `a`.

Atomicity

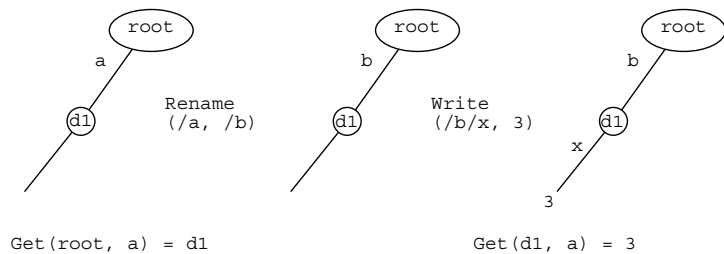
The `MemNames` and `ObjNames0` specs made all the update operations atomic. For code to satisfy these specs, it must hold some kind of lock on every directory touched by `GetDN`, or at least on the name looked up in each such directory. This can involve a lot of directories, and since the name space is a graph it also introduces the danger of deadlock. It’s therefore common for systems to satisfy only the weaker atomicity spec of `ObjNames`, which says that looking up a simple name is atomic, but the entire lookup process is not.

This means that `Read(/a/x)` can return `3` even though there was never any instant at which the path name `/a/x` had the value `3`, or indeed was defined at all. To see how this can happen, suppose:

```
initially /a is the directory d1 and /b is undefined;
initially x is undefined in d1;
concurrently with Read(/a/x) we do Rename(/a, /b); Write(/b/x, 3).
```

The following sequence of actions yields `Read(/a/x) = 3`:

```
In the Read, Get(root, a) = d1
Rename(/a, /b) makes /a undefined and d1 the value of /b
Write(/b/x, 3) makes 3 the value of x in d1
In the Read, RET d1.get(x) returns 3.
```

Obviously, whether this possibility is important or not depends on how clients are using the name space.

Coherence

Other things being equal, everyone prefers a coherent or ‘sequentially consistent’ memory, in which there is a single order of all the concurrent operations with the property that the result of every read is the result that a simple memory would return after it has done all the preceding writes in order. Maintaining coherence has costs, however, in the amount of synchronization that is required if parts of the memory are cached, or in the amount of availability if the memory is replicated. We will discuss the first issue in detail at the end of the course. Here we consider the availability of a replicated memory.

Recall the majority register from the beginning of the course. It writes a majority of the replicas and reads from a majority, thus ensuring that every read must see the most recent write. However, this means that you can’t do either a read or a write unless you can talk to a majority. There we used a general notion of majority in which the only requirement is that every two majorities have a non-empty intersection. Applying this idea, we can define separate read and write *quorums*, with the property that every read quorum intersects every write quorum. Then we can make reads more available by making every replica a read quorum, at the price of having the only write quorum be the set of all replicas, so that we have to do every write to all the replicas.

An alternative approach is to weaken the spec so that it’s possible for a read to see old values. We have seen a version of this already in connection with crashes and write buffering, where it was possible for the system to revert to an old state after a crash. Now we propose to make the spec even more non-deterministic: you can read an old value at any time, and the only restriction is that you won’t read a value older than the most recent *sync*. In return, we can now have much more availability in the code, since both a read and a write can be done to a single replica. This means that if you do `write(/a, 3)` and immediately read `a`, you may not get `3` because the `read` might use a different replica that hasn’t seen the `write` yet. Only *sync* requires communication among the replicas.

We give the spec for this as a variation on `ObjNames`. We allow `nil` to be in `dd(n)`, representing the fact that `n` has been undefined in `dd`.

```

TYPE DD          = N -> SEQ Z                               % remember old values
APROC GetFromS(d, n) -> Z = <<                             % we write d.get(n)
% The non-determinism wouldn't be allowed if this were a function
VAR z | z IN s(d)(n) => RET z [*] RET nil >>               % return any old value

PROC SetToS(d, n, z) =                                     % we write d.set(n, z)
  s(d)(n) := [(s(d)!n => s(d)(n) [*] {} ) + {z}]           % add z to the state

PROC Sync(pn) RAISES {error} =
  VAR d, n, z |
    (d, n) := GetDN(pn, true); z := s(d)(n).last;
    IF z # nil => s(d)(n) := {z} [*] s(d) := s(d){n -> } FI

```

This spec is common in the naming service for a distributed system, for instance in the Internet’s DNS or Microsoft’s Active Directory. The name space changes slowly, it isn’t critical to see the very latest value, and it is critical to have high availability. In particular, it’s critical to be able to look up names even when network partitions make some working replicas unreachable.

Schemas

In the database world, a schema is the definition of what names are defined (and usually also of the type of each name’s value).⁵ Network management calls this a ‘management information base’ or MIB. Depending on the application there are very different rules about how the schema is defined.

In a file system, for example, there is usually no official schema written down. Nonetheless, each operating system has conventions that in practice have the force of law. A Unix system without `/bin` and `/etc` will not get very far. But other parts of the name space, especially in users’ private directories, are completely variable.

By contrast, a database system takes the schema very seriously, and a management system takes at least some parts of it seriously. The choice has mainly to do with whether it is people or programs that are using the name space. Programs tend to be much less flexible; it’s a lot of work to make them adapt to missing data or pay attention to unexpected additional data

Minor issues

We mention in passing some other, less fundamental, ways in which the specs for name spaces differ.

Rules about overwriting. Some systems allow any name to be overwritten, others treat directories, or non-empty directories, specially to reduce the consequences of careless errors.

Access control. Many systems enforce rules about which users or programs are allowed to read or write various parts of the name space.

⁵ Gray and Reuter, *Transaction Processing*, Morgan Kaufmann, 1993, pp 768-786.

Resource control. Writes often consume resources that are expensive or in fixed supply, such as disk blocks. This means that they can fail if the resources are exhausted, and there may also be a quota system that limits the resource consumption of users or programs.

Roots

It's not turtles all the way down.

Anonymous

So far we have ducked the question of how the `root` is represented, or the `D` in a link that plays a similar role. In `ObjNames0` we said `D = Int`, leaving its interpretation entirely to the `s` component of the state. In `ObjNames` we said `D` is a pair of procedures, begging the question of how the procedures are represented. The representation of a root depends entirely on the implementation. In a file system, for instance, a root names a disk, a disk partition, a volume, a file system exported from a server, or something like that. Thus there is another name space for the roots (another level of indirection). It works in a wide variety of ways. For example:

In MS-DOS, you name a physically connected disk drive. If the drive has removable media and you insert the wrong one, too bad.

On the Macintosh, you use the string name of a disk. If the system doesn't know where to find this disk, it asks the user. If you give the same name to two removable disks, too bad.

On Digital VMS, disks have unique identifiers that are used much like the string names on the Macintosh.

For the NFS network file system, a root is named by a host name or IP address, plus a file system name or handle on that host. If that name or address gets assigned to another machine, too bad.

In a network directory a root is named by a unique identifier. There is also a set of servers that might store replicas of that directory.

In the secure file system, a root is named by the hash of a public encryption key. There's also a network address to help you find the file system, but that's only a hint.⁶

In general it is a good idea to have absolute names (unique identifiers) for directories. This at least ensures that you won't use the wrong directory if the information about where to find it turns out to be wrong. A UID doesn't give much help in locating a directory, however. The possibilities are:

Store a set of places to look along with the UID. The problem is keeping this set up to date.

Keep another name space that maps UID's to locations (yet another level of indirection).

The problem is keeping this name space up to date, and making it sufficiently available.

For the former, every location can register itself periodically. For the latter, replication is good. We will talk about replication in detail later in the course.

Search some ad-hoc set of places in the hope of finding a copy. This search is often called a 'broadcast'.

⁶ Mazières, Kaminsky, Kaashoek, and Witchel, Separating key management from file system security. *Proc. 17th ACM Symposium on Operating Systems Principles*, Dec. 1999. www.pdos.lcs.mit.edu/papers/sfs:sosp99.pdf.

We defined the interface routines to start from a fixed `root`. Some systems, such as Unix, have provisions for changing the root; the `chroot` system call does this for a process. In addition, it is common to have a more local context (called a 'working directory' for a file system), and to have syntax to specify whether to start from the root or the working directory (presence or absence of an initial '/' for a Unix file system).

Examples

These are to expand your mind and to help you recognize a name space when you come across it under some disguise.

File system directory Example: `/udir/lampson/pocs/handouts/12-naming`

Not a tree, because of `.` and `..`, hard links, and soft links.

Devices, named pipes, and other things can appear as well as files.

Links and mounting are important for assembling the name space you want.

Files may have attributes, which are a little directory attached to the file.

Sometimes resources, fonts, and other OS rigmarole are stored this way.

inodes There is a single inode directory, usually coded as a function rather than a table: you compute the location of the inode on the disk from the number.

For system-wide inodes, prefix a system-wide file system or volume name.

Plan⁷ This operating system puts all its objects into a single name space: files, devices, pipes, processes, display servers, and search paths (as union directories).

Semantic file system⁸ Not restricted to relational databases.

Free-text indexing: `~lampson/Mail/inbox/ (word="compiler")`

Program cross-reference: `/project/sources/ (calls="DeleteFile")`

Table (relational data base)	Example:	<i>ID no (key)</i>	<i>Name</i>	<i>Salary</i>	<i>Married?</i>
		1432	Smith	21,000	Yes
		44563	Jones	35,000	No
		8456	Brown	17,000	Yes

We can view this as a naming tree in several ways:

`#44563/Name = Jones` key's value is a `D` that defines `Name`, `Salary`, etc.

`Name/#44563 = Jones` key's value is the `Name` field of its row

The second way, `cat Name/*` yields

Smith Jones Brown

⁷ Pike et al., The use of name spaces in Plan 9, *ACM Operating Systems Review* **27**, 2, Apr. 1993, pp 72-76.

⁸ Gifford et al., Semantic file systems, *Proc. 13th ACM Symposium on Operating System Principles*, Oct. 1991, pp 16-25 (handout 13).

Network naming ⁹	<p>Example: <code>theory.lcs.mit.edu</code></p> <p>Distributed code. Can share responsibility for following the path between client and server in many ways.</p> <p>A directory handle is a machine address (interpreted by some communication network), plus some id for the directory on that machine.</p> <p>Attractive as top levels of complete naming hierarchy.</p>
E-mail addresses	<p>Example: <code>rinard@lcs.mit.edu</code></p> <p>This syntax patches together the network name space and the user name space of a single host. Often there are links (called forwarding) and directories full of links (called distribution lists).</p>
SNMP ¹⁰	<p>Example: Router with circuits, packets in circuits, headers in packets, etc.</p> <p>Internet Simple Network Management Protocol</p> <p>Roughly, view the state of the managed entity as a table, treating it as a name space the way we did earlier. You can read or write table entries.</p> <p>The <code>Next</code> action allows a client to explore the name space, whose structure is read-only. Ad hoc <code>write</code> actions are sometimes used to modify the structure, for instance by adding a row to a table.</p>
Page tables	<p>Divide up the virtual address, using the first chunk to index a first level page table, later chunks for lower level tables, and the last chunk for the byte in the page.</p>
I/O device addressing	<p>Example: Memory bus. SCSI controller, by device register addresses. SCSI device, by device number 0..7 on SCSI bus. Disk sector, by disk address on unit.</p> <p>Usually there is a pure read/write interface to the part of the I/O system that is named by memory addresses (the device registers in the example), and a message interface to the rest (the disk in the example).</p>
Multiplexing a channel	<p>Examples: Node-node network channel $\rightarrow n$ process-process channels. Process-kernel channel $\rightarrow n$ inter-process channels. ATM virtual path $\rightarrow n$ virtual circuits.</p> <p>Given a channel, you can multiplex it to get sub-channels. Sub-channels are identified by addresses in messages on the main channel. This idea can be applied recursively, as in all good name spaces.</p>
LAN addresses	<p>48-bit ethernet address. This is flat: the address is just a UID.</p>

⁹ B. Lampson, Designing a global name service, *Proc. 4th ACM Symposium on Principles of Distributed Computing*, Minaki, Ontario, 1986, pp 1-10. RFC 1034/5 for DNS.

¹⁰ M. Rose, *The Simple Book*, Prentice-Hall, 1990.

Hierarchical network addresses ¹¹	<p>Example: <code>16.24.116.42</code> (an IP address).</p> <p>An address in a big network is hierarchical.</p> <p>A router knows its parents and children, like a file directory, and also its siblings (because the parent might be missing)</p> <p>To route, traverse up the name space to least common ancestor of current place and destination, then down to destination.</p>
Network reference ¹²	<p>Example: <code>6.24.116.42/11234/1223:44 9 Jan 1995/item 21</code></p> <p>Network address + port or process id + incarnation + more multiplexing + address or export index.</p> <p>Some applications are remote procedure binding, network pointer, network object</p>
Abbreviations	<p>A, talking to B, wants to pass a big value v, say a font or security credentials.</p> <p>A makes up a short name n for v (sometimes called a ‘cookie’, though it’s not the same as a Web cookie) and passes that.</p> <p>If B doesn’t know n’s value v, it calls back to A to get it, and caches the result.</p> <p>Sometimes A tells v to B when it chooses n, and B is expected to remember it.</p> <p>This is not as good because B might run out of space or fail and restart.</p>
World Wide Web	<p>Example: <code>http://ds.internic.net/ds/rfc-index.html</code></p> <p>This is the URL (Uniform Resource Locator) for Internet RFCs.</p> <p>The Web has a read/write interface.</p>
Spec names	<p>Example: <code>ObjNames.Enum</code></p>
Telephone numbers	<p>Example: <code>1-617-253-6182</code></p>
Postal addresses	<p>Example: <code>Prof. Butler Lampson Room 43-535 MIT Cambridge, MA 02139</code></p>

¹¹ R. Perlman, *Connections*, Prentice-Hall, 1993.

¹² Andrew Birrell et al., Network objects, *Proc. 14th ACM Symposium on Operating Systems Principles*, Asheville, NC, Dec. 1993 (handout 25).