# 21
# User-Defined Types and Procedural Data Structures as Complementary Approaches to Data Abstraction

## John C. Reynolds[1]
*Syracuse University*

**Abstract** *User-defined types (or modes) and procedural (or functional) data structures are complementary methods for data abstraction, each providing a capability lacked by the other. With user-defined types, all information about the representation of a particular kind of data is centralized in a type definition and hidden from the rest of the program. With procedural data structures, each part of the program which creates data can specify its own representation, independently of any representations used elsewhere for the same kind of data. However, this decentralization of the description of data is achieved at the cost of prohibiting primitive operations from accessing the representations of more than one data item. The contrast between these approaches is illustrated by a simple example.*

## Introduction

User-defined types and procedural data structures have both been proposed as methods for data abstraction, i.e., for limiting and segregating the portion of a program which depends upon the representation used for some kind of data. In this paper we suggest, by means of a simple example, that these methods are complementary, each providing a capability lacked by the other.

The idea of user-defined types has been developed [Morris 73,74; Liskov 74; Fischer 73; Wulf 77] and has its roots in earlier work [Dahl 72b*]. In this approach, each particular conceptual kind of data is called a *type*, and for each type used in a program, the program is divided into two parts: a type definition and an "outer" or "abstract" program. The type definition specifies the representation to be used for the data type and a set of primitive operations (and perhaps constants), each defined in terms of

the representation. The choice of representation is hidden from the outer program by requiring all manipulations of the data type in the outer program to be expressed in terms of the primitive operations. The heart of the matter is that any consistent change in the data representation can be effected by altering the type definition without changing the outer program.

Various notions of procedural (or functional) data structures have been developed [Reynolds 70*; Landin 65; Balzer 67*]. In this approach, the abstract form of data is characterized by the primitive operations which can be performed upon it, and an item of data is simply a procedure or collection of procedures for performing these operations. The essence of the idea is seen most clearly in its implementation: an item of procedural data is a kind of record called a *closure* which contains both an internal representation of the data and a pointer (or flag field) to code for procedures for manipulating this representation. A program with access to a closure record is only permitted to examine or access the internal representation by executing the code indicated by the pointer, so that this code serves to close off or protect the internal representation.

In comparison with user-defined types, procedural data structures provide a decentralized form of data abstraction. Each part of the program which creates procedural data will specify its own form of representation, independently of the representations used elsewhere for the same kind of data, and will provide versions of the primitive operations (the components of the procedural data item) suitable for this representation. There need be no part of the program, corresponding to a type definition, in which all forms of representation for the same kind of data are known. But a price must be paid for this decentralization: a primitive operation can have access to the representation of only a single data item, namely the item of which the operation is a component.

Apparently this price is inevitable. If an operation is to have access to the representation of more than one item of data, each of which may have several possible representations, then its definition cannot be "decentralized" into one part for each representation, since one must provide for every possible *combination* of representations. Presumably this requires the definition to occur at a point in the program where all possible representations of the operands are known.

## Linguistic preliminaries

Before illustrating these ideas, we must digress to explain (informally) the language we will use. It is an applicative language, similar to pure LISP [McCarthy 60] or the applicative subsets of GEDANKEN [Reynolds 70*], PAL [Evans 68], or ISWIM [Landin 66], but with a complete type structure somewhat like ALGOL 68 [Van Wijngaarden 75]. Types will be indicated

by writing $\in \tau$, where $\tau$ is a type expression, after binding occurrences of identifiers (except where the type is obvious from context). Type expressions are constructed with the operators $\rightarrow$ denoting functional procedures, $\times$ denoting a Cartesian product, and $+$ denoting a named disjoint union.

The named disjoint union is sufficiently novel to require a more detailed explanation. If $\tau_1, \ldots, \tau_n$ are type expressions denoting the sets $S_1, \ldots, S_n$ and $i_1, \ldots, i_n$ are distinct identifiers, then

$$i_1 : \tau_1 + \cdots + i_n : \tau_n$$

is a type expression denoting the set of pairs

$$\{\langle i_k, x \rangle \mid 1 \leqslant k \leqslant n \text{ and } x \in S_k\}.$$

If $e$ is an expression of type $\tau_k$ with value $x$, then

$$\textbf{inject } i_k \, e$$

is an expression of type $i_1 : \tau_1 + \cdots + i_n : \tau_n$ with value $\langle i_k, x \rangle$.

Let $e$ be an expression of type $i_1 : \tau_1 + \cdots + i_n : \tau_n$ with value $\langle i, x \rangle$, let $i_{k_1}, \ldots, i_{k_m}$ be distinct members of the set of identifiers $\{i_1, \ldots, i_n\}$, for $1 \leqslant j \leqslant m$ let $l_j$ be an expression of type $\tau_{k_j} \rightarrow \tau'$ with value $f_j$, and let $e'$ be an expression of type $\tau'$ with value $x'$. Then

$$\textbf{unioncase } e \textbf{ of } \left( i_{k_1} : l_1, \ldots, i_{k_m} : l_m, \textbf{other} : e' \right)$$

is an expression of type $\tau'$ with the value

$$\textbf{if} \begin{bmatrix} i = i_{k_1} \\ \vdots \\ i = i_{k_m} \\ \text{otherwise} \end{bmatrix} \textbf{then} \begin{bmatrix} f_1(x) \\ \vdots \\ f_m(x) \\ x' \end{bmatrix}$$

When $m = n$, the **other** clause will be omitted.

We use the type expression *nilset* to denote a standard one-element set, whose unique member is denoted by { }.

# Integer sets as a user-defined type

Our example is an implementation of the abstract concept of sets of integers. Using the approach of user-defined types, we wish to define a type *set* and primitive constants and functions

$$none \in set$$

$$all \in set$$

$$limit \in \textbf{integer} \times \textbf{integer} \times set \rightarrow set$$

$$union \in set \times set \rightarrow set$$

$$exists \in \textbf{integer} \times \textbf{integer} \times set \rightarrow \textbf{Boolean}$$

satisfying the specifications

$$none = \{\ \}$$

$all = $ The set of all (machine-representable) integers

$$limit(m,n,s) = s \cap \{k | m \leqslant k \leqslant n\}$$

$$union(s1,s2) = s1 \cup s2$$

when $m \leqslant n$, $exists(m,n,s) = (\exists k) m \leqslant k \leqslant n$ and $k \in s$

To make our solution seem more realistic, we require that the execution of *limit* and *union* should require time and space bounded by constants which are independent of their arguments. Of course this will exact a price in the speed of *exists*.

An appropriate and simple solution is to represent a set by a list structure which records the way in which the set is constructed via primitive operations. Thus the representation of a set is a disjoint union, over the four set-valued primitive functions (including constants), of sets of possible arguments for these functions. More precisely, this representation is defined by the recursive type declaration:

$$set = nonef : \textbf{nilset} + allf : \textbf{nilset} + limitf : \textbf{integer} \times \textbf{integer} \times set$$
$$+ unionf : set \times set$$

and the effect of *none, all, limit,* or *union* is to imbed its arguments into the appropriate kind of list element:

$$none = \textbf{inject}\ nonef\ (\ )$$

$$all = \textbf{inject}\ allf\ (\ )$$

$$limit(m,n,s) = \textbf{inject}\ limitf(m,n,s)$$

$$union(s1,s2) = \textbf{inject}\ unionf(s1,s2)$$

(Roughly speaking, we are representing sets by a free algebra with constants *none* and *all*, and operators *limit* and *union*.) The entire computational burden of interpreting this representation falls upon the function *exists*:

```
exists(m,n,s) = unioncase s of
    (nonef : λ( ). false,
    allf : λ( ). true,
    limitf : λ(m1,n1,s1). max(m,m1) ⩽ min(n,n1)
    and exists(max(m,m1),min(n,n1),s1),
    unionf : λ(s1,s2). exists(m,n,s1) or exists(m,n,s2))
```

(We assume that the operations **and** and **or** do not evaluate their second operand when the first operand is sufficient to determine their result.)

Although the above is a definition of the type *set* which meets our specifications, it can be easily improved, even within the time and space

constraints imposed upon *limit* and *union*. For example, both *limit* and *union* can be optimized by taking advantage of some obvious properties of sets—the result of *limit* can be simplified when its last argument is *none* or another application of *limit*, and the result of *union* can be simplified when either argument is *none* or *all*:

*limit* $(m,n,s)$ = **unioncase** *s* **of**
    $(nonef:\lambda(\ ).$ *none*,
    $limitf:\lambda(m1,n1,s1).$ **if** $max(m,m1)\leqslant min(n,n1)$
            **then inject** $limitf(max(m,m1),\ min(n,n1),\ s1)$ **else** *none*,
    **other** : **inject** $limitf\ (m,n,s))$

*union* $(s1,s2)$ = **unioncase** *s*1 **of**
    $(nonef:\lambda(\ ).$ *s*2,
    $allf:\lambda(\ ).$ *all*,
    **other** :**unioncase** *s*2 **of**
            $(nonef:\lambda(\ ).$ *s*1,$allf:\lambda(\ ).$ *all*,
            **other** : **inject** $unionf\ (s1,s2)))$

In conclusion, we show how our specification of integer sets might be "packaged" in a language permitting user-defined types:

**newtype** *set* = *nonef* : **nilset** + *allf* : **nilset** + *limitf* : **integer** × **integer** × *set*
    + *unionf* : *set* × *set*
**with** *none* ∈ *set* = **inject** *nonef* ( ),
    *all* ∈ *set* = **inject** *allf* ( ),
    *limit* ∈ **integer** × **integer** × *set* → *set* =
        $\lambda(m,n,s).$ **unioncase** *s* **of**
        $(nonef:\lambda(\ ).$ *none*,
        $limitf:\lambda(m1,n1,s1).$ **if** $max(m,m1)\leqslant min(n,n1)$
            **then inject** $limitf\ (max(m,m1),\ min(n,n1),\ s1)$ **else** *none*,
        **other** : **inject** $limitf\ (m,n,s))$,
    *union* ∈ *set* × *set* → *set* =
        $\lambda(s1,s2).$ **unioncase** *s*1 **of**
        $(nonef:\lambda(\ ).$ *s*2, $allf:\lambda(\ ).$ *all*,
        **other** : **unioncase** *s*2 **of**
            $(nonef:\lambda(\ ).$ *s*1, $allf:\lambda(\ ).$ *all*,
            **other** : **inject** $unionf\ (s1,s2)))$,
    *exists* ∈ **integer** × **integer** × *set* → **Boolean** =
        $\lambda(m,n,s).$ **unioncase** *s* **of**
        $(nonef:\lambda(\ ).$ **false**,
        $allf:\lambda(\ ).$ **true**,
        $limitf:\lambda(m1,n1,s1).$ $max(m,m1)\leqslant min(n,n1)$
                **and** $exists(max(m,m1),min(n,n1),s1)$,
        $unionf:\lambda(s1,s2).$ $exists(m,n,s1)$ **or** $exists(m,n,s2))$
**in** $\langle$*outer program*$\rangle$

The language used here is an outgrowth of the ideas discussed in [Reynolds 74*]. A complete exposition of this language is beyond the scope of this paper, but the following salient points should be noted.

(1) The type declaration between **newtype** and **with** binds all occurrences of the type identifier *set* throughout the above expression (including occurrences in ⟨*outer program*⟩). The ordinary declarations between **with** and **in** bind all occurrences of the ordinary identifiers *none, all, limit, union,* and *exists* throughout the expression.

(2) With regard to occurrences of *set* between **with** and **in**, the type declaration behaves like a mode definition in ALGOL 68, i.e., *set* is equivalent to the type expression on the right side of the type declaration, and the type-correctness of the text in **with...in** depends upon this type expression.

(3) In ⟨*outer program*⟩, occurrences of *set* behave like a primitive type, e.g., **integer** or **Boolean**. In other words, ⟨*outer program*⟩ must be a correctly typed expression regardless of what type expression might be equivalent to *set*. This insures that all manipulations of the user-defined type in ⟨*outer program*⟩ must be expressed in terms of the primitives declared in **with...in**.

(4) Although it is not illustrated by our example, it should be possible to declare simultaneously several related user-defined types between **newtype** and **with**. This ability is needed to permit the definition of multiargument primitive functions which act upon more than one user-defined type. An example might be the use of the types *point* and *line* in a program for performing geometrical calculations.

# Integer sets as procedural data structures

We now develop integer sets as procedural data structures. The starting point is the realization that all we ever want to do to a set $s$, aside from using it to construct other sets, is to evaluate the Boolean expression $exists(m,n,s)$. This suggests that we can simply equate the set $s$ with the Boolean function $\lambda(m,n).\ exists(m,n,s)$ that characterizes the only information we want to extract from the set.

Thus we define

$$set = \textbf{integer} \times \textbf{integer} \rightarrow \textbf{Boolean}$$

and specify that if $s \in set$ represents the "mathematical" set $s_0$, then for $m \leqslant n$,

$$s(m,n) = (\exists k) m \leqslant k \leqslant n \text{ and } k \in s_0.$$

The need for defining the primitive function *exists* has vanished since this function has been *internalized*—its value for a particular *set* is simply the (only component of the) *set* itself. The remaining primitive constants

and functions are easily defined by:

$none = \lambda(m,n).$ **false**
$all = \lambda(m,n).$ **true**
$limit(m,n,s) = \lambda(m1,n1).$
    $max(m,m1) \leqslant min(n,n1)$ **and** $s(max(m,m1),min(n,n1))$
$union(s1,s2) = \lambda(m,n).$ $s1(m,n)$ **or** $s2(m,n)$

In this approach, there is no $\langle$*outer program*$\rangle$ from which the definition $set = $ **integer** $\times$ **integer** $\rightarrow$ **Boolean** is hidden. Any part of the program can create a *set* by giving an appropriate function whose internal representation (the collection of values of global variables which form the fields of the closure record) can be arbitrary. For example, in augmenting an existing program, one might write

$$\lambda(m,n). \; even(m) \text{ or } (m < n)$$

to denote the set of even integers, or

**letrec** $s = \lambda(m,n). \; (m \leqslant n) \text{ **and** } (p(m) \text{ **or** } s(m+1,n)) \text{ **in** } s$

to denote the set of integers satisfying the predicate $p$. The procedural approach insures that these definitions will mesh correctly with the rest of the program, even though they introduce novel representations.

This kind of extensional capability, which is the main advantage of the procedural approach, is offset by two limitations. In the first place, although (ignoring computability considerations) every set can be represented by a function in **integer** $\times$ **integer** $\rightarrow$ **Boolean**, the converse is false. To represent a set, a function $s$ must satisfy

$$s(m,n) = \bigvee_{k=m}^{n} s(k,k)$$

for all $m$ and $n$ such that $m \leqslant n$. This kind of condition, which cannot be checked syntactically, must be satisfied by all parts of the program which create sets.

A more important limitation is that only the function *exists*, which has been internalized as (the only component of) a procedural data item, is truly primitive in the sense of having access to the internal representation of a set. Essentially, we have been forced to express the functions *limit* and *union* in terms of the internalized *exists*. We are fortunate that our example permits us to do this at all. Even so, we are prevented from optimizing *limit* and *union* as we did in the user-defined-type development. There is no practically effective way that $limit(m,n,s)$ can "see" whether $s$ has the form *none* or $limit(m1,n1,s1)$, or that $union(s1,s2)$ can "see" whether $s1$ or $s2$ has the form *none* or *all*.

In fact, this difficulty can be surmounted for *limit* but not for *union*. The solution is to internalize *limit* as well as *exists*, so that both functions have access to internal representations. Thus we represent sets by pairs of

functions:
$$set = (integer \times integer \rightarrow Boolean) \times (integer \times integer \rightarrow set)$$
and specify that if $s$ represents the mathematical set $s_0$ then for $m \leqslant n$,
$$s.1(m,n) = (\exists k)m \leqslant k \leqslant n \text{ and } k \in s_0,$$
and for all $m$ and $n$, $s.2(m,n)$ represents the mathematical set
$$s_0 \cap \{k | m \leqslant k \leqslant n\}$$
(Here $s.1$ and $s.2$ denote the components of the pair $s$.)

In this approach, we may define *none* by:
$$none = (\lambda(m,n). \textbf{ false}, \lambda(m,n). \text{ none})$$
Note the peculiar kind of recursion which is characteristic of this style of programming: the second component of *none* is a function which does not call itself but rather returns itself as a component of its result.

To define *all* and *union* we first define an "external" *limit* $\in$ **integer** $\times$ **integer** $\times$ *set* $\rightarrow$ *set* which will be called upon by the internal limiting functions (i.e., the second components) of *all* and *union*:

$limit(m,n,s) =$
    $(\lambda(m1,n1).max(m,m1) \leqslant min(n,n1) \textbf{ and } s.1(max(m,m1),$
        $min(n,n1)),$
    $\lambda(m1,n1).\textbf{if } max(m,m1) \leqslant min(n,n1) \textbf{ then}$
        $limit(max(m,m1), min(n,n1), s) \textbf{ else } none)$

Then

    $all = (\lambda(m,n). \textbf{ true}, \lambda(m,n). \text{ limit}(m,n,all))$
    $union(s1,s2) = (\lambda(m,n). \text{ s1}.1(m,n) \textbf{ or } s2.1(m,n),$
        $\lambda(m,n). \text{ limit}(m,n,union(s1,s2))).$

With these definitions, the internal limiting functions perform simplifications analogous to those performed by *limit* in the user-defined-type approach. Indeed, if one examines the behavior of the closures which would represent sets in an implementation of this definition, one finds that they mimic the list structures of the type approach almost exactly (except for the simplifications performed by union).

But even to someone who is experienced with procedural data structures, the internalization of *limit* is more a tour de force than a specimen of clear programming. Moreover, internalization cannot be applied to give a function such as *union* access to the internal representation of more than one argument, i.e., we could convert *union*(s1,s2) to a component of s1 or of s2 but not both.

# Conclusions

In comparison with user-defined types, procedural data structures offer a more decentralized method of data abstraction which precludes any interaction between different representations of the same kind of data. This

offers the advantage of easier extensibility at the price of prohibiting primitive operations from accessing the representations of more than one data item.

Of course, the two approaches can be combined. For example, we can augment our user-defined-type definition to include an additional primitive *functset* $\in$ (integer $\times$ integer$\rightarrow$**Boolean**)$\rightarrow$*set* which accepts a functional set (in the sense of the first part of the previous section) and produces an equivalent value of type *set*. It is sufficient to add one more kind of record to the disjoint union defining *set* and one more alternative to the branches defining *exists*:

**newtype** *set* = $\cdots$ + *functsetf* : (integer $\times$ integer$\rightarrow$**Boolean**)
**with** ...
    *functset* $\in$ (integer $\times$ integer$\rightarrow$**Boolean**)$\rightarrow$*set* = $\lambda f.$ **inject** *functsetf f*,
    *exists* $\in$ integer $\times$ integer $\times$ *set*$\rightarrow$**Boolean** =
        $\lambda(m,n,s).$ **unioncase** *s* **of**
        ( ... *functsetf* : $\lambda f.$ $f(m,n)$)
**in** $\langle$*outer program*$\rangle$

However, this kind of combination is hardly a unification. To some extent, the data-representation structuring approach of [Dahl 72b*] unifies the concepts of user-defined types and procedural data structures, but only at the expense of combining their limitations. It appears that this is inevitable, that the two concepts are inherently distinct and complementary.

The reader should be cautioned that this is a working paper describing ongoing research. In particular, the linguistic constructs we have used are tentative and will require considerable study and evolution before they can be integrated into a complete programming language. The extension of these constructs to languages with imperative features is a particularly murky area.