

TOPICS IN PROGRAMME SPECIFICATION AND DESIGN:

SPECIFICATION AND DESIGN
OF
DISTRIBUTED SYSTEMS

I.H. SORENSEN
WOLFSON COLLEGE

A THESIS SUBMITTED FOR THE DEGREE OF DOCTOR OF PHILOSOPHY
IN THE UNIVERSITY OF OXFORD, SEPTEMBER 1981.

TOPICS IN PROGRAMME SPECIFICATION AND DESIGN :
SPECIFICATION AND DESIGN OF
DISTRIBUTED SYSTEMS

Ib Holm Sorensen
Wolfson College

A thesis submitted for the degree of Doctor of Philosophy in the University of Oxford,
September 1981.

ABSTRACT.

This thesis presents a method for *specifying, analysing and refining* the designs of *distributed systems*. Distributed systems are systems which consist of several autonomous process components.

The characteristics of the specification method employed in this thesis can be summarised as follows - 1) The structure of a system's specification indicates the structure of the system's realisation, 2) A design is specified entirely in terms of the permissible activity across the interfaces between process components (i.e. the communications); such a specification gives the rules for the behaviour of each process component and postpones decisions about its internal structure, 3) permissible activity is described in terms of predicates on the history of past communications.

This specification method will be shown to allow important questions about the behaviour of a distributed system to be posed early in the design process; in particular designs will be analysed with respect to termination and absence of deadlocks.

The specification method can be employed to describe systems in different degrees of detail, and it is demonstrated that a specification can evolve to a stage close to realisation using a stepwise refinement method which ensures that the important properties are maintained.

Table of Contents.

1. INTRODUCTION	1
1.1. THE OUTLINE OF THIS THESIS	2
2. THE SPECIFICATION LANGUAGE	5
2.1. THE SET THEORETICAL BASIS	6
2.1.1. Sets and Operations on Sets	6
2.1.2. Relations and Operations on Relations	8
2.1.3. Functions	9
2.2. THE SCHEMA NOTATION	11
2.2.1. Examples of the Use of Schemas	11
2.2.2. Schema Renaming	14
2.2.3. Generic Schemas	16
2.3. EXTENSIBILITY	16
2.3.1. Operations on Relations and Functions	16
2.3.2. Operations on Sequences	18
2.4. Notational Conventions	20
3. NON DETERMINISTIC SYSTEMS	21
3.1. FORMALISATION OF NON DETERMINISTIC BEHAVIOUR	22
3.2. A BOUNDED BUFFER	29
3.3. FORMALISATION OF NON DETERMINISTIC SYSTEMS	35
3.3.1. Deadlock and the Invariant Condition	37
3.3.2. Termination and the Variant Function	39
4. SPECIFICATION OF DISTRIBUTED SYSTEMS	45
4.1. MOTIVATION FOR DISTRIBUTION	45
4.2. THE INDEPENDENT MODULE	46
4.2.1. The Classical Independent Module	45
4.2.2. The Process Module	45
4.2.3. The Readiness Principle	48
4.3. DISTRIBUTED SPECIFICATIONS AND DISTRIBUTED REALISATIONS	50
4.4. ENVIRONMENTS FOR SYSTEMS	55
4.4.1. A Requesting Environment	55
4.4.2. A Selfcommitting Environment	55
4.4.3. An Interactive Environment	56
4.5. ANALYSIS OF THE BEHAVIOUR OF DISTRIBUTED SYSTEMS	57
5. EXAMPLES OF SPECIFICATIONS	61
5.1. A BOUNDED BUFFER	61
5.2. MISCELLANEOUS SYSTEMS	65
5.2.1. Buffers	65
5.2.2. Queues	66
5.3. CRITICAL SECTIONS AND MUTUAL EXCLUSION	68
5.4. PREVENTION OF DEADLOCK	72

5.5. A COMMUNICATION NETWORK	78
5.5.1. A Discussion of the Design of a Station	78
5.5.2. A Specification of a Network	84
5.5.3. An Analysis of the Network	85
6. DECOMPOSITION	89
6.1. INTRODUCTORY EXAMPLE	90
6.2. VALID DECOMPOSITION	92
6.3. LIMITATIONS OF THE METHOD	96
6.3.1. Piped Systems	96
6.3.2. Unordered Pipes	97
6.3.3. A Resource Monitor	99
7. CONCLUSION	101
7.1. RELATED WORK	102
7.2. FUTURE RESEARCH	103
LIST OF REFERENCES	105

Acknowledgement.

I am indebted to my colleagues at the University of Aarhus, Denmark, in particular Peter Kornerup and Michael Spler for supporting my application for the grant which enabled me to join the software engineering research team at Oxford University. I am also grateful to Professor Tony Hoare for accepting me as a D.Phil. student in his department.

I would like to thank the staff at the Programming Research Group, Oxford University for their valuable guidance in research and general assistance during the preparation of this thesis. Thanks are especially given to Jean Raymond Abrial, Tim Clement, Bernard Sufrin and Cliff Jones.

It should be mentioned that the work presented in this thesis is a consequence of Jean Raymond Abrial's pioneering research in the area of formal specifications. His insight and suggestions have inspired many of the ideas presented in this thesis.

Special thanks are due to Hilary Darragh for enriching my English language and to Tim Clement and Steve Schuman for their constructive criticisms while editing my thesis.

lb Holm Sorensen

1. INTRODUCTION.

The tools and methods used in the creation of computerised systems are being intensively studied by many researchers with the aim of improving our capability for constructing 'good' computer systems.

The quality of such systems depends to a large extent on the methods used for documentation during the development process.

The *documentation* of a system has three functions: firstly, it identifies the *needs* which the system is intended to fulfil; secondly, it specifies the decisions taken by the designers of the system as to *how* those needs are to be fulfilled; and thirdly, it gives a detailed description of how the design decisions are realised. The documentation of the needs must be *precise* in order for it to serve as a basis for a contract between the customer and the supplier. The documentation of the design and the realisation must be precise in order to *test* and *verify* whether a suggested implementation fulfills that contract. Furthermore, since the *needs* will change over time it is important that all parts of the documentation (including the programs which document the realisation) are easy to modify.

These concerns have inspired the development of numerous methodologies for the creation of software. Current techniques use *informal languages* or *diagrams* as description tools for documenting the requirements, the design decisions and the structure of the realisation. The disadvantage of such techniques lies in the ambiguity inherent in the notational framework. Among other things, this ambiguity makes it impossible to *formally test* a design against a requirement specification, or a realisation against the design decisions (*program verification*). *Theoretical* or *formal* tools which do not have these drawbacks have been proposed, but system developers are reluctant to use such formal methods because they are claimed to be inadequate for the documentation of *real systems* - a claim which is understandable since their usefulness has yet to be demonstrated on large commercial products.

But, it is incontestable that formal methods are necessary, since only the use of unambiguous formal descriptions enables us to *verify* whether decisions taken in the development process are consistent with the original specification of the needs. Formal methods for developing systems which consist of modules running in isolation and having well-defined starting and stopping points are fairly well understood [Jones,14]. These methods, in which modules can be described as mathematical functions, cannot, however, be employed in the development of systems which consist of several interacting autonomous process-components (i.e. *distributed systems*). The 'function' computed by a subcomponent in such a system can at any point in time be *interrupted* and influenced by non-deterministically occurring stimuli from the environment or from other concurrently executing process-components. The component itself may interfere with computations within other components depending on some non local information.

The complexity of such systems necessitates the development of alternative formal methods. The aim of this thesis is to contribute to this development by providing tools for specifying, analysing and refining the designs of distributed systems.

1.1. THE OUTLINE OF THIS THESIS.

CHAPTER 2 gives an informal presentation of the *language* used in this thesis. The language is an extension of a conventional set theoretical notation. The extensions are simple syntactical conventions and shorthands, which are introduced in order to improve the writability and readability of specifications. The language has been developed independently of the special area - specification of distributed systems - to which it is being applied in this thesis. The main aim in the development of this new language has been to provide a uniform notation and a formal framework for reasoning and proving properties about computer systems. The language (in its current form) has been shown to meet this aim (e.g. [Sufrin.21]). It is not within the scope of this thesis to give a formal definition for the syntax and semantics of the notation in question. Hence, chapter 2 gives only an informal overview of the language, and is included mainly to ease the understanding of the specifications in the remaining parts of the thesis.

Specifications of a system include a description of an abstract state-space and a description of the transformations of that state. The state is constrained by axioms. These constraints proscribe those 'abstract' transformations which would violate the axioms, and allow (by definition) those which would not. In other words, an operational model in which the transformations may or may not occur whenever sufficient pre-conditions are met can be derived from the description of the abstract state and the constraints. The approach, for specifying distributed systems, presented in this thesis allows for systems behaviour to be investigated using such an operational model for describing *non-deterministic systems*. A formal framework for discussing *non-deterministic behaviour* is developed in CHAPTER 3. Although the work presented in that chapter is original, it is not entirely novel as it embodies a recasting of familiar concepts (e.g. well-founded relations [Abrial.1] and correctness of non-deterministic constructs [Dijkstra.7]). The theory of non-deterministic systems developed in chapter 3 is generally applicable to the analysis of systems whose behaviour can be characterised by a relation.

CHAPTER 4 introduces a new approach for describing *distributed systems*. Such systems are described in terms of constraints on the communications which may occur across a set of interfaces (i.e. connections). The constraints are imposed by axioms involving the *past history of communications* along one or more of the connections. A communication history along a single connection is modelled as a *sequence*. The specification of a distributed system provides enough information for a non-deterministic model of its behavior to be defined and analysed (using the framework developed in chapter 3). Furthermore, the inside of each process-component (an object between a set of interfaces) can be constructed from the the axioms which constrain the activities along its own *external interfaces*.

In CHAPTER 5 the specification method is applied to a number of small systems. Questions which are normally of concern when designing distributed systems (mutual exclusion, starvation, deadlock etc.) are discussed in connection with these examples.

The specification method can be used to describe distributed systems at different levels of abstraction. In a very abstract form only a few interfaces are visible - and in the least abstract form all interfaces from the realisation are visible. Hence, it can support a *stepwise development* method. Chapter 6 gives a method for comparing (the behaviour of) different levels of description, thus giving us the basis for *refining* specifications.

CHAPTER 7 presents some conclusions and proposes topics for further study.

2.1. THE SET THEORETICAL BASIS.

This section is an overview of the basic notation of the proposed language. No attempt is made to define the semantics of the constructs formally; this and succeeding subsections are intended only to provide enough information about the language to permit discussion about specifications written therein.

2.1.1. Sets and Operations on Sets.

The following expression denotes a particular subset of a set X , namely the set of elements for which PRED holds,

$$\{ x:X \mid \text{PRED}(x) \} \quad \text{DEF (2.1)}$$

where $\text{PRED}(x)$ stands for a predicate in which x appears free. This is a rewriting of the classical form,

$$\{ x \mid x \in X \wedge \text{PRED}(x) \}$$

The set above is *non-empty* if and only if

$$(\exists x : X \mid \text{PRED}(x))$$

The explicit construction of a finite set is denoted as usual by enumerating its elements, as in

$$\{ a, b, c, d \}$$

The following expressions denote, respectively, the set of all subsets of the set X and the set of all finite subsets of the set X

$$\mathcal{P}(X), \quad \mathcal{F}(X)$$

The empty set is denoted by

$$\{ \}$$

The number of elements in a finite set X is

$$\text{card}(X)$$

We shall use the standard set operations of union, intersection and difference as well as the inclusion and membership operators. They are denoted as usual by

$$\cup \quad \cap \quad - \quad \subset \quad \not\subset \quad \subseteq \quad \in \quad \notin$$

Furthermore

$$\cup \quad \cap$$

will be used as prefix operators to denote the distributed union and intersection respectively. The use of these operators take a special form, e.g. let $\text{EXP}(i)$, for all i in an index set I , denote a subset of X , then the following

$$\cup \{ \text{EXP}(i) \mid i : I \}$$

is a rewriting of

$$\cup \{ s : \mathcal{P}(X) \mid (\exists i : I)(s = \text{EXP}(i)) \}$$

The *introduction* (or the *type declaration*) of a variable x , which can only be bound to values from the set X , is written,

$$x : X$$

(see the examples above)

A particular element with the value EXP from

$$\{ x:X \mid \text{PRED}(x) \}$$

can be written

$$(\mu x : X \mid \text{PRED}(x)) (x = \text{EXP}) \quad \text{DEF (2.2)}$$

provided, of course, that $\text{PRED}(\text{EXP})$ holds.

The designation of an *arbitrary* element from the same set is written

$$(\mu x : X \mid \text{PRED}(x))$$

provided that

$$\{ x : X \mid \text{PRED}(x) \} \neq \{ \}$$

Construction of a subset of the Cartesian Product

$$X \times Y$$

is written

$$\{ x:X; y:Y \mid \text{PRED}(x,y) \} \quad \text{DEF (2.3)}$$

Instead of

$$\{ (x,y) \mid (x,y) \in X \times Y \wedge \text{PRED}(x,y) \}$$

2.1.2. Relations and Operations on Relations.

The set of all binary relations from a set X to a set Y is denoted by

$$X \leftrightarrow Y$$

which simply is a shorthand for

$$\mathcal{P}(X \times Y)$$

A finite relation can be constructed explicitly as follows.

$$\{ x_1 \leftrightarrow y_1, x_2 \leftrightarrow y_2, x_1 \leftrightarrow y_2 \}$$

which maps x_1 to y_1 and y_2 etc.

Given a binary relation R from X to Y , then the *image* of a subset SX of X through R i.e.

$$\{ y : Y \mid (\exists x : SX) (x,y) \in R \}$$

is denoted by

$$R(SX)$$

The *inverse* of a binary relation R from X to Y is denoted by

$$R^{-1}$$

Given a relation $R : X \leftrightarrow Y$ we have

$$R^{-1} \in Y \leftrightarrow X$$

and

$$\begin{aligned} \text{dom}(R) &= R^{-1}(Y) \\ \text{ran}(R) &= R(X) \end{aligned}$$

which denote the *domain* and the *range* of the relation R .

When expressing the relationship between two elements we use the well-known form

$$x R y.$$

which is syntactically equivalent to

$$y \in R(\{x\})$$

The *composition* of two binary relations

$$\begin{aligned} R_1 &: X \leftrightarrow Y \\ R_2 &: Y \leftrightarrow Z \end{aligned}$$

is denoted by

$$R_2 \circ R_1$$

which is a relation from X to Z for which

$$\begin{aligned} (\forall x:X; z:Z) \\ x (R_2 \circ R_1) z \iff (\exists y:Y) (x R_1 y \wedge y R_2 z) \end{aligned}$$

2.1.3. Functions.

The set of *partial* functions from X to Y is a subset of the set of all relations between X and Y , and is denoted by

$$X \rightarrow Y$$

$$\begin{aligned} X \rightarrow Y = \\ \{ r:X \leftrightarrow Y \mid (\forall y:Y \mid y \in \text{ran}(r)) r(r^{-1}(\{y\})) = \{y\} \} \end{aligned}$$

The construction of a partial function 'f' from X to Y , which maps an element in X to an element of Y given by the expression $\text{EXP}(x)$, might be written.

$$f = (\lambda x : X \mid \text{PRED}(x)) (\text{EXP}(x)) \quad \text{DEF (2.4)}$$

Provided that

$$(\forall x:X \mid \text{PRED}(x)) (\text{EXP}(x) \in Y)$$

we have

$$\text{dom}(f) = \{ x:X \mid \text{PRED}(x) \}$$

Given two predicates PRE and POST then we may construct a function h of type

$$h : \{x:X \mid \text{PRE}(x)\} \rightarrow Y$$

as follows

$$h = (\lambda x:X | PRE(x)) ((\mu y:Y | POST(x,y)) (y = EXP(x))) \quad \text{DEF (2.5)}$$

We have

$$\begin{aligned} \text{dom}(h) &\subseteq \{x:X | PRE(x)\} \\ \text{ran}(h) &\subseteq \{y:Y | (\exists x:X | POST(x,y))\} \end{aligned}$$

The domain of h can be determined to be

$$\begin{aligned} \text{dom}(h) &= \{x:X | PRE(x)\} \cap \\ &\quad \{x:X | (\exists y:Y | POST(x,y)) \wedge y = EXP(x)\} \end{aligned}$$

The set of *total functions* from X to Y is a subset of all partial functions over these sets, and is denoted by

$$X \rightarrow Y$$

$$X \rightarrow Y = \{ f:X \rightarrow Y \mid \text{dom}(f) = X \}$$

Consider the definition of h above: if

$$(\forall x:X | PRE(x)) POST(x, EXP(x))$$

then h is a total function i.e.

$$h \in \{x:X | PRE(x)\} \rightarrow Y$$

hence the 'transformations' denoted by

$$y = EXP(x)$$

can be considered as a realisation of the specification

$$h_{\text{spec}} = (\lambda x:X | PRE(x)) ((\mu y:Y) (POST(x,y)))$$

(cf. [Jones.14])

Given $f : X \rightarrow Y$, then the single element in the *image* of a single element through a function is written as usual,

$$f(x)$$

where

$$f(x) = (\mu y : Y) (y \in f(\{x\}))$$

We can define a function which maps a *non-empty* set into an *arbitrary element* of that set as follows

$$\begin{aligned} \tau &: \mathcal{P}(X) \rightarrow X; \\ \tau &= (\lambda s:\mathcal{P}(X) | s \neq \{\}) ((\mu x:X) (x \in s)) \end{aligned}$$

2.2. THE SCHEMA NOTATION.

The previous section provided a uniform notation for the implicit construction of sets, elements, relations and functions (see DEF(2.1-5))

The *syntactic* schema

$$a: A; b: B \dots z: Z \mid \text{predicate}(a,b,\dots,z)$$

appears in the characterisation of sets and of the domain and range of functions.

This schema is also used in the quantified expressions of the Predicate Calculus, e.g. the expressions (appearing in section 2.1.1. and 2.1.3)

$$\begin{aligned} &(\exists x:X | PRED(x)) \\ &(\forall x:X | PRED(x)) (EXP(x) \in Y) \end{aligned}$$

are rewritings of

$$(\exists x) (x \in X \wedge PRED(x)),$$

and

$$(\forall x) ((x \in X \wedge PRED(x)) \Rightarrow (EXP(x) \in Y))$$

The notation introduced in this thesis allows the *definition* and *naming* of such syntactic schemas.

2.2.1. Example of the Usage of Schemas.

Let the Natural Numbers be denoted by

$$\mathbf{N}$$

We can give the name TWONUM to the schema

$$n:\mathbf{N}; m:\mathbf{N} \mid n > m$$

either by writing

$$\text{TWONUM} \equiv [n:\mathbf{N}; m:\mathbf{N} \mid n > m]$$

or by using a *vertical* presentation.

TWONUM

$n : \mathbf{N};$
$m : \mathbf{N}$
$n > m$

The part above the horizontal line within the box is referred to as the *signature*, the part below as the *axiom* or the *invariant*.

The schema TWONUM can be used as a *textual macro* in

- 1) The implicit construction of a set.

twonum = {TWONUM}
i.e.
twonum = {n:N; m:N | n ≥ m}

- 2) The implicit construction of a function.

subtract : {TWONUM} → N
subtract = (λ TWONUM) (n-m)

- 3) The construction of elements.

pair = (μ TWONUM) (n=3; m=2)

NB. the definition is 'valid', since 3 ≥ 2

- 4) In quantified expressions.

(∃ TWONUM) (n = m+10)
(∀ TWONUM) (n + 1 > m)

- 5) The definition of *new* schemas, by *extending* existing schemas

INTERVAL
TWONUM; set : P(N)
(∀ i : s) (m ≤ i ∧ i ≤ n)

which is a shorthand for

INTERVAL
n, m : N; set : P(N)
n > m; (∀ i : s) (m ≤ i ∧ i ≤ n)

- 6) Definition of *new* schemas by *combining* existing schemas.

Let

LIMSET
set : P(N);
card(set) < 256

then we might define

SET2
TWONUM; LIMSET
n = card(set); m = card({ j : set j ≠ 0 })

Note that the combination of two schemas may introduce 'duplicated' declarations in the signature of the resulting schema. Such a schema is identical to a schema in which one of the declarations is removed, e.g. the schema

[INTERVAL; LIMSET]

in which the name set 'occurs' twice is

[INTERVAL; card(set) < 256]

- 7) In *Theorems*, as illustrated below.

The signature part of a schema may be empty, hence we can define.

SUM_and_AVG
n = x + y; m = (x + y)/2

NB '/' is integer division, i.e.

/ : N × N → N
(∀ i, j, r : N) (r=i/j ⇔ 0 ≤ i - m*j < j)

Given the schema,

PARAM

$$x : \mathbf{N}; y : \mathbf{N1}$$

NB, $\mathbf{N1}$ denotes the *non-zero* natural numbers

we can define a function which gives the *sum* and the *average* of two numbers in the following way.

$$\text{sumavg} = (\lambda \text{ PARAM}) (\mu \text{ TWONUM}) (\text{SUM_and_AVG})$$

which defines a total function if we can prove the following theorem.

$$x \in \mathbf{N}; y \in \mathbf{N1}; n = x+y; m = (x+y)/2$$

$$\vdash$$

$$n \in \mathbf{N}; m \in \mathbf{N}; n \geq m$$

which leads to the last important use of the schema notation. The theorem above can be written as follows.

$$\text{PARAM}; \text{SUM_and_AVG} \vdash \text{TWONUM}$$

2.2.2. Schema Renaming.

The bound variables of a schema (i.e. the variables which appear in the signature part of a schema) may be systematically *renamed*, e.g.

$$\text{TWONUM}'$$

denotes the text

$$n' : \mathbf{N}; m' : \mathbf{N} \mid n' > m'$$

and

$$\text{TWONUM}_{\text{ex}}$$

denotes

$$n_{\text{ex}} : \mathbf{N}; m_{\text{ex}} : \mathbf{N} \mid n_{\text{ex}} > m_{\text{ex}}$$

The schema notation allows individual renaming of one or several variables of the signature, e.g.

$$\text{TWONUM}[k/n]$$

denotes

$$k : \mathbf{N}; m : \mathbf{N} \mid k \geq m$$

It is important to understand, however, that

$$\{\text{TWONUM}\} = \{\text{TWONUM}'\} = \{\text{TWONUM}_{\text{ex}}\} = \{\text{TWONUM}[k/n]\}$$

i.e. systematic renaming of bound variables of a schema does not affect the set which corresponds to that schema, cf.

$$(\lambda x)(x+2) = (\lambda y)(y+2)$$

Additionally, the schema notation allows for renaming of the free variables in a schema, i.e. given

$$\text{A_NUM} \equiv [n:\mathbf{N} \mid n < \text{Limit}]$$

then

$$\text{A_NUM1} \equiv \text{A_NUM}[L1/\text{Limit}]$$

denotes a new schema. Note that, if $L1 \neq \text{Limit}$ then

$$\{\text{A_NUM}\} \neq \{\text{A_NUM1}\}$$

The main use of schema renaming is in the definition of state transformations, where, informally speaking, the same state variables appear twice.

Example:

let us define a simple state as follows.

$$\text{STATE}$$

$$\text{counter} : \mathbf{N}$$

A decrement function,

$$\text{decr} : \{\text{STATE}\} \rightarrow \{\text{STATE}\}$$

can be defined using the renaming facility

$$\text{decr} = (\lambda \text{ STATE}) (\mu \text{ STATE}') (\text{counter}' = \text{counter} - 1)$$

NB

$$\text{dom}(\text{decr}) = \{\text{STATE} \mid \text{counter} > 0\}$$

2.2.3. Generic Schemas.

Informal statements or definitions, such as those in section 1 of this paper, often take a *generic* form. Consider the following definition which is generic with respect to X and Y

"let X and Y be sets. The *partial injections* from X to Y are exactly the partial functions from X to Y whose *inverses* are also partial functions."

In the schema notation this concept is formalised, e.g. the above statement will have the form

par_inj	X	Y
$f : X \rightarrow Y$		
$f^{-1} \in Y \rightarrow X$		

A schema may have several *instances*, e.g. the set of partial injections over the natural numbers is

{ par_inj[N,N] }

and the set of partial injections from N to the Cartesian Product over N is,

{ par_inj[N,N×N] }

2.3. Extensibility.

The importance of the possibility to extend the formal framework in which one wants to present a specification has already been emphasised in the introduction.

The notation provided for extending a theory allows for the introduction of new objects, which are either constructively or axiomatically defined, and for the introduction of new operations on these objects. Such extensions can be grouped and the groups are given names, which enable the writers of a specification to reference theories developed elsewhere.

2.3.1. Operations on Relations and Functions.

The elementary Set Theory which was presented in section 2.1.1 can be extended to include a set of useful operations on relations. Let us first introduce the *identity* function (relation) over a subset of any set, i.e.

REL-OP	X
$id : P(X) \rightarrow (X \leftrightarrow X)$	
$(\forall s : P(X))$ $id(s) = \{ (x,x) \mid x \in s \}$	

This definition is *generic* in X , hence, we have several instances of the identity function. The use of different instances of a definition is illustrated in the next extension which is *generic* with respect to S and T . This extension defines two *domain* restriction operators and a *codomain* restriction operator.

REL-OP1	S	T
$op(\uparrow) : (S \leftrightarrow T) \times P(S) \rightarrow S \leftrightarrow T ;$		
$op(\downarrow) : (S \leftrightarrow T) \times P(T) \rightarrow S \leftrightarrow T ;$		
$op(\setminus) : (S \leftrightarrow T) \times P(S) \rightarrow S \leftrightarrow T$		
$(\forall r : S \leftrightarrow T; s : P(S); t : P(T))$ $(r \uparrow s = r \circ id[S](s) \wedge$ $ r \downarrow t = id[T](t) \circ r \wedge$ $ r \setminus s = r \uparrow (S-s))$		

As for the identity function we have several instances of the domain-restriction operators, e.g.

$\uparrow [N,N] \in (N \leftrightarrow N) \times P(N) \rightarrow N \leftrightarrow N$
 $\uparrow [N,P(N)] \in (N \leftrightarrow P(N)) \times P(N) \rightarrow N \leftrightarrow P(N)$

The 'parameter-list' of the instantiation of an operator is often omitted, since its 'type' can be determined from the context in which the operator appears, e.g.

let f, g, S be declared as follows

$f, g : N \rightarrow N ;$
 $S : P(N)$

then

$g = f \uparrow S$

means

$g = f \uparrow [N,N] S$

The n th iterate of a relation

$R : S \leftrightarrow S$

can informally be written as

$$R^n = \underbrace{R \circ R \circ \dots \circ R}_{n\text{-times}}$$

The following extensions define the *iteration* operators.

REL-OP2	X
$op(\wedge) : (X \leftrightarrow X) \times \mathbf{N} \rightarrow (X \leftrightarrow X) ;$	
$op(\ast) : (X \leftrightarrow X) \rightarrow (X \leftrightarrow X)$	
$(\forall r : X \leftrightarrow X; n : \mathbf{N}1)$	
$(r^0 = id(X) \quad \wedge$	
$r^n = r^{(n-1)} \circ r)$	
$(\forall r : X \leftrightarrow X) r^\ast = \cup \{r^n \mid n : \mathbf{N}\}$	

NB. the *infix* operator \wedge is implied (as usual) by writing the second operand as a superscript to the first operand.

In the following we add the function *overriding* operator to the theory of functions.

FUNC-OP	S T
$op(\oplus) : ((S \rightarrow T) \times (S \rightarrow T)) \rightarrow (S \rightarrow T)$	
$(\forall g, f : S \rightarrow T) (g \oplus f = (g \setminus \text{dom}(f)) \cup f)$	

NB. functions are sets to which the Set operations can be applied.

2.3.2. Operations on Sequences.

A theory which is often used in the specification of systems is the theory of sequences. In the following we will formalise *sequences* and some operations on them.

Let us first formalise the well known 2-dots, as in

3..m

$op(\dots) : \mathbf{N} \times \mathbf{N} \rightarrow \mathcal{P}(\mathbf{N})$
$(\forall n, m, p : \mathbf{N}) (p \in n..m \iff n \leq p \wedge p \leq m)$

Sequences are then formalised as a subset of the partial functions, e.g. a *sequence* of length 4 over a set

{a, b, c, d, f, g}

might be defined as.

{(1,a), (2,c), (3,b), (4,a)}

This sequence will be denoted by

< a c b a >

The empty sequence is denoted by

<>

Sequences are 'formally' added to our theory by the following generic extension.

SEQ	X
$seq, seq1 : \mathcal{P}(\mathbf{N} \rightarrow X)$	
$(\forall s : seq)$	
$(\text{dom}(s) \in \mathcal{F}(\mathbf{N}) \wedge$	
$\text{dom}(s) = 1..card(\text{dom}(s)))$	
$seq1 = seq - \{\langle \rangle\}$	

In the following we add some operators to our 'current theory'

SEQ-OP1	X
$\# : seq[X] \rightarrow \mathbf{N};$	
$next : seq[X] \rightarrow (X \leftrightarrow X);$	
$first, last : seq1[X] \rightarrow X;$	
$ending,$	
$beginning : seq1[X] \rightarrow seq[X]$	
$(\forall s : seq)$	
$(\#(s) = card(s) \quad \wedge$	
$next(s) = s \circ succ \circ s^{-1})$	
$(\forall s : seq1)$	
$(first(s) = s(1) \quad \wedge$	
$last(s) = s(\#(s)) \quad \wedge$	
$ending(s) = s \circ (succ \setminus \{0\}) \quad \wedge$	
$beginning(s) = s \uparrow (1.. \#(s) - 1)$	

The function *succ* is the successor function over the Natural Numbers.

The concatenation operator for sequences is defined by.

SEQ-OP2	X
$op(*) : seq[X] \times seq[X] \rightarrow seq[X]$	
$(\forall s, s' : seq) s * s' = s \cup (s' \circ pred^{*(s)})$	

The function *pred* is the predecessor function over the Natural Numbers.

2.4. NOTATIONAL CONVENTIONS.

The *schema* notation permits the grouping of concepts into named structures. The names are used in the construction of sets, functions etc.. We also allow for the *naming of theory-extensions* (as in SEQ and SEQ_OP1 from section 2.3.2) simply in order to be able to reference specific theories elsewhere in the document. The names are also suggestive, i.e. SEQ, SEQ_OP1 and SEQ_OP2 make up the current theory of SEQUENCES. Even with such a naming facility the cross-referencing in a large document (like a thesis) becomes unmanageable, therefore this thesis makes use of an indexing system, in which all important sections of formal text have been indexed. *Schemas* and *Extensions* are categorised as definitions, hence we use DEF as a qualifier for them (see DEF(3.2)). *Theorems*, *Lemmas* and *Figures* are also isolated from the informal text using 'boxes'; their index-numbers are prefixed with the qualifiers THEOREM, LEMMA and FIGURE respectively. Theorems are when convenient given the name of the *theory* to which they belong, see LOOP:DEF(3.2) and the theorems LOOP:THEOREM(3.3-6).

3. NON-DETERMINISTIC SYSTEMS.

The purpose of this chapter is to investigate the properties of non-deterministic behaviour. The results will be used to justify the introduction of a particular approach to the specification of distributed systems.

The descriptions of most systems do not include complete descriptions of control mechanisms. This means that we cannot fully predict the behaviour of the system by looking at its description - the behaviour of our system is *non-deterministic*. The *top-down* development process which takes a system through design stages to a realisation, involves, among other things, adding further decisions about actual behaviour to the description; hence the non-determinism present at the top level might well be resolved at some lower level where a more detailed description of the behaviour is given.

We can therefore expect, even when developing sequential deterministic programs, to be faced with the problem of proving properties of designs which involve non-deterministic behaviour. Systems which are eventually implemented as strictly sequential programs may at a certain stage in the design process best be described using non-deterministic constructs [Dijkstra,8]. The reasons for this are twofold - 1) the designer should not impose in advance any decisions about the system which may have an adverse effect on future refinements and realisations - 2) the properties of a design are often independent of decisions about actual dynamic behaviour, hence a description of these decisions is pure 'noise' and should be left out.

Furthermore, there is a very important class of systems which cannot be realised using sequential control constructs because their actual behaviour is controlled by or synchronized with *external events* (i.e. events activated by the environment) which occur non-deterministically. For example, we expect the control-flow of the program for an air-traffic control system to be influenced by the unpredictable event of an aeroplane entering into the airspace monitored by the system; we expect a terminal multiplexer to acknowledge a key-stroke; we expect the activity of a batch operating system to be interrupted upon entry of a new deck of cards.

System behaviour may also include *internal events* that occur non-deterministically, i.e. may or may not occur whenever the preconditions which determine their possible occurrence are satisfied, as for a communication between components in a multi-process system. Common to the realisations of these systems is the fact that they are composed of several highly independent processing components, whose progress between points of synchronization is unknown (distributed concurrent systems). Such systems cannot be realised using purely sequential control mechanisms, since any decision that attempts to prescribe the order of occurrence for 'inter-process' events in advance (in either the design or implementation) could easily cause unnecessary delays at run-time, and even infinite delays (i.e. deadlocks).

In the first section of this chapter a simple model for non-deterministic behaviour will be introduced. Definitions of a set of functions, which will be used to characterise non-deterministic behaviour, will be given. This model and these definitions are applied to a very simple example in the section which follows. The last section will summarise the properties of non-deterministic behaviour to be considered, and give a method for proving those properties.

3.1. FORMALISATION OF NON-DETERMINISTIC BEHAVIOUR.

Non-deterministic behaviour may be pictured as a *state diagram* (a directed graph) in which the *nodes* represent the *states* of the system and the outgoing edges represent the *events*. A node with two or more outgoing edges represents a state of the system where sufficient conditions for the occurrence of more than one event are satisfied. Nodes with no outgoing edges represent states where no further progress can be made - these states are referred to as *Terminal-states*.

Even for small systems this pictorial method for describing non-deterministic behaviour becomes unmanageable, and can therefore only be used to illustrate special dynamic properties - never complete systems. A more suitable method of description is needed.

The events (the *state-transitions*) will be described as *partial functions*. Whenever the system which is being described is in a state belonging to the domain of one of these partial functions, the corresponding event may occur, thus transforming the system. Consequently in such a model the total behaviour is represented by a *relation* (usually referred to as the *Step* - relation) and the *Terminal* - states are states which do not belong to the domain of this relation.

We can define a *schema* for describing such a model :

ND	S	DEF (3.1)
Events	$\mathcal{P}(S \rightarrow S)$;	
Step	$S \leftrightarrow S$;	
Terminal	$\mathcal{P}(S)$	
Step	$= \cup(\text{Events})$;	
Terminal	$= S - \text{dom}(\text{Step})$	

Note that only systems with a finite number of events are considered; most of the results in this chapter do however apply to the more general but unrealistic case where

$$\text{Events} : \mathcal{P}(S \rightarrow S)$$

For any given state s , $\text{Step}(\{s\})$ is the set of states which can be reached from s in one step, i.e. if s is in the domain of two events

(ev1, ev2)
then
 $\text{Step}(\{s\}) = \{\text{ev1}(s), \text{ev2}(s)\}$

The relation

Step^n

will, applied to state s , describe the states which can be reached after n steps.

The domain of the relation,

$\text{Step}^n \downarrow \text{Terminal}$

is the set of states in S in which the system may be started and possibly terminate after n steps. This intuition justifies the following definition:

LOOP	S	DEF (3.2)
loop	$(S \leftrightarrow S) \rightarrow (S \leftrightarrow S)$	
	$(\forall r : S \leftrightarrow S)$ $(\text{loop}(r) = \cup\{r^n \downarrow (S - \text{dom}(r)) \mid n : \mathbb{N}\})$	

The operator *loop* can be used to establish *partial correctness*. (A system is *partially correct* if it produces the desired result or fails to terminate).

If the system is started (from s) within the domain of *loop* for a relation *Step* then the system may terminate, and furthermore if it terminates the resulting state is within $\text{loop}(\text{Step})(\{s\})$. This may be expressed as follows:

LOOP	S	THEOREM (3.3)
R	$S \leftrightarrow S$;	
s	S ;	
$s \in \text{dom}(\text{loop}(R))$	\vdash	
$(\exists n : \mathbb{N})(\exists s' : S)((s R^n s') \wedge s' \notin \text{dom}(R))$;		
$(\forall n : \mathbb{N}; s' : S)$		
$(s R^n s' \wedge s' \notin \text{dom}(R)) \Rightarrow s' \in \text{loop}(R)(\{s\})$		

The theorem below is a reformulation of the definition of *loop*, a reformulation which follows directly from the definition of the *transitive closure* (R^*) for relations and the distributivity of \downarrow over \cup .

LOOP $\frac{}{S}$ THEOREM (3.4)
 $R : S \leftrightarrow S \vdash$
 $\text{loop}(R) = R^* \downarrow (S - \text{dom}(R))$

A system started in a Terminal-state will remain in that state; in other words states belonging to the complement of the domain of a relation are in the domain of the loop of the same relation. This is a consequence of -

LOOP $\frac{}{S}$ THEOREM (3.5)
 $R : S \leftrightarrow S \vdash$
 $\text{loop}(R) = \text{id}(S - \text{dom}(R)) \cup (\text{loop}(R) \bullet R)$

A rule of iteration for non-deterministic systems may be formulated as follows:

LOOP $\frac{}{S}$ THEOREM (3.6)
 $R : S \leftrightarrow S;$
 $\text{INV} : \mathcal{P}(S)$
 $R(\text{INV}) \subseteq \text{INV} \vdash \text{loop}(R)(\text{INV}) = \text{INV} - \text{dom}(R)$

proof :

- (1) $R(\text{INV}) \subseteq \text{INV}$ (hypothesis)
- (2) $R^*(\text{INV}) \subseteq \text{INV}$ (by induction, as follows)
 - (a) $R^0(\text{INV}) \subseteq \text{INV}$ (the basis of the induction)
 - (b) $R^n(\text{INV}) \subseteq \text{INV}$ (induction hypothesis)
 - (c) $R(R^n(\text{INV})) \subseteq R(\text{INV})$ (b)
 - (d) $R^{n+1}(\text{INV}) \subseteq \text{INV}$ (c), (1)
- (3) $(R^* \downarrow (S - \text{dom}(R)))(\text{INV}) \subseteq \text{INV}$ ((2) and def \downarrow)
- (4) $(R^* \downarrow (S - \text{dom}(R)))(\text{INV}) \subseteq (S - \text{dom}(R))$ (def \downarrow)
- (5) $(\text{loop}(R)(\text{INV})) \subseteq (\text{INV} \cap (S - \text{dom}(R)))$ (3,4, TH3.4)
- (6) $\text{id}(S - \text{dom}(R)) \subseteq \text{loop}(R)$ (TH3.5)
- (7) $\text{id}(S - \text{dom}(R))(\text{INV}) \subseteq \text{loop}(R)(\text{INV})$ (6)
- (8) $(\text{INV} - \text{dom}(R)) \subseteq (\text{loop}(R)(\text{INV}))$ (7)
- (9) $(\text{loop}(R)(\text{INV})) = (\text{INV} - \text{dom}(R))$ (5), (8)

Intuition -

Given a state S . Suppose INV_CON is an invariant condition for a system which behaves according to a relation Step and suppose we can describe $S - \text{dom}(\text{Step})$ as $\{ S \mid \text{TERM_CON} \}$. The invariant condition holds after any number of iterations (occurrences of events) provided it was satisfied upon initialization, hence this condition also holds when (IF!) the system terminates; furthermore we know that no more progress is possible i.e. the condition

$$\text{INV_CON} \wedge \text{TERM_CON}$$

holds, i.e. the system terminates in a state belonging to

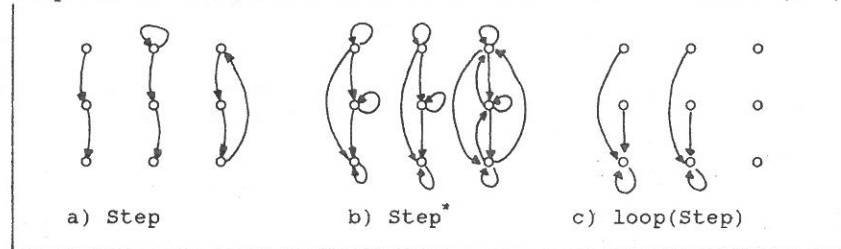
$$\{ S \mid \text{INV_CON} \} - \text{dom}(\text{Step})$$

The following graphs (FIGURE 3.7) illustrate the behaviour of a simple system.

We observe the following :

- 1) A system having the transition-relation Step will never terminate if it is started outside the domain of $\text{loop}(\text{Step})$.
- 2) A system started within that domain may terminate, but termination is not guaranteed.

Loop and Termination FIGURE (3.7)



To be able to distinguish between properties such as possible termination and guaranteed termination we introduce -

CYCLE $\frac{}{S}$ DEF (3.8)

$\text{Span} : (S \leftrightarrow S) \rightarrow \mathcal{P}(S \rightarrow S);$
 $\text{Cycle} : (S \leftrightarrow S) \rightarrow (S \leftrightarrow S)$

$(\forall r : S \leftrightarrow S)$
 $(\text{Span}(r) = \{f : S \rightarrow S \mid \text{dom}(f) = \text{dom}(r) \wedge f \subseteq r\}) \wedge$
 $\text{Cycle}(r) = \text{loop}(r) \uparrow \cap \{\text{dom}(\text{loop}(sp)) \mid sp : \text{Span}(r)\}$

NB, the set of spanning functions for a relation r is non-empty, as

$$g = (\lambda x : S \mid x \in \text{dom}(r))(\tau(r(\{x\})))$$

$$g \in \text{Span}(r)$$

A spanning-function, sp , for a relation R

$$sp \in \text{Span}(R)$$

characterises a particular deterministic behaviour obtained by making decisions in advance about the 'future' non-deterministic behaviour. In other words a particular spanning-function together with a starting state describes a run-time trace.

The set of spanning-functions for a system contains functions which describe 'worst-case' behaviour, e.g. infinite loops. Hence if we start a system in a state from which an infinite loop can be reached then there exists a spanning-function.

say sp , such that the particular starting state is outside the domain of $loop(sp)$. On the other hand if a system is started within the domains of the loops of all spanning-functions, then termination is guaranteed. (as an example, see FIGURE(3.27) and FIGURE(3.28) which illustrate a relation and two of the spanning-functions for that relation). Using these definitions we can give a more precise characterization of non-deterministic behaviour.

A rewriting gives -

$$\text{CYCLE} \frac{S}{\vdash (\forall r : S \leftrightarrow S) \quad \text{THEOREM (3.9)}} \\ \text{Cycle}(r) = \cup \{ loop(sp) \mid sp : \text{Span}(r) \} \uparrow \\ \quad \quad \quad \cap \{ \text{dom}(loop(sp)) \mid sp : \text{Span}(r) \}$$

The next theorem characterises the deterministic behaviour of a system if it is started in a terminal-state (outside the domain of the step- relation).

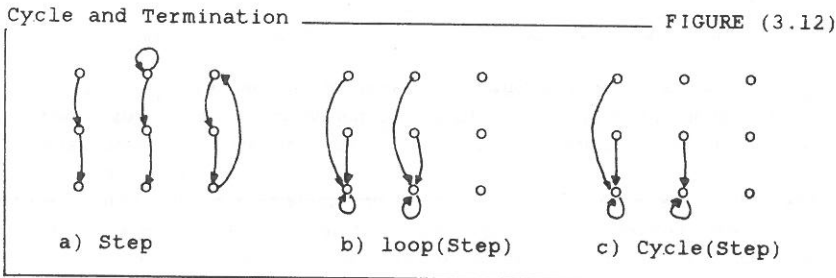
$$\text{CYCLE} \frac{S}{\vdash (\forall r : S \leftrightarrow S) \quad \text{THEOREM (3.10)}} \\ \text{Cycle}(r) = \text{id}(S - \text{dom}(r)) \cup \\ ((loop(r) \circ r) \uparrow \cap \{ \text{dom}(loop(sp)) \mid sp : \text{Span}(r) \})$$

proof :
use DEF(3.8) and THEOREM(3.5) to rewrite $loop(r)$.

The rule of iteration for non-deterministic systems can now be strengthened-

$$\text{CYCLE} \frac{S}{\begin{array}{l} r : S \leftrightarrow S \\ \text{INV} : \mathcal{P}(S) \\ r(\text{INV}) \subseteq \text{INV} \end{array} \vdash \text{Cycle}(r)(\text{INV}) = (\text{INV} - \text{dom}(r)) \quad \text{THEOREM (3.11)}}$$

Furthermore, as mentioned earlier, the domain of $Cycle$ gives the set of states for a system from which it is guaranteed to terminate; this is illustrated by the graphs in FIGURE(3.12) which is a continuation of the description of the small system from FIGURE(3.7) -



The following Termination theorem expresses that there is, for each state within the domain of the $Cycle$ of a finite branching relation (R), i.e. a relation for which-

$$(\exists n : \mathbb{N}) (\forall x : \text{dom}(R)) (\text{card}(R(x)) < n)$$

an upper limit to the number of transitions which can be performed before a terminal-state is reached. A proof of the theorem can be given using the Rank function, which will be defined in DEF(3.14).

$$\text{TERMINATION} \frac{S}{\begin{array}{l} \text{Events} : \mathcal{F}(S \rightarrow S); \\ s : S; \\ s \in \text{dom}(\text{Cycle}(\cup \text{Events})) \quad \vdash \\ (\exists \text{Max} : \mathbb{N}) \\ (\forall n : \mathbb{N}; s' : S) ((s (\cup \text{Events})^n s') \Rightarrow (n < \text{Max})) \end{array} \quad \text{THEOREM (3.13)}}$$

The Rank function defines the maximum distance to a terminal state - i.e. the maximum distance to the leaves of the relation in question. The definition of Rank is justified by definition DEF(3.15), which defines the distance-set, and THEOREM(3.16) in which it is proved that all distance-sets of finite branching relations are finite sets. Hence, a distance-set has a maximum.

$$\text{RANK} \frac{S}{\begin{array}{l} \text{Rank} : \mathcal{F}(S \rightarrow S) \rightarrow (S \rightarrow \mathbb{N}) \\ \\ (\forall \text{events} : \mathcal{F}(S \rightarrow S)) \\ (\forall s : S \mid s \in \text{dom}(\text{Cycle}(\cup \text{events}))) \\ \text{Rank}(\text{events})(s) = \max(\text{distance-set}(\cup \text{events}, s)) \end{array} \quad \text{DEF (3.14)}}$$

For a relation R , the distance-set is defined only for elements which are in the domain of $Cycle(R)$ -

$$\text{DISTANCES} \frac{S}{\begin{array}{l} \text{distance-set} : (S \leftrightarrow S) \times S \rightarrow \mathcal{P}(\mathbb{N}) \\ \\ (\forall R : S \leftrightarrow S; s : S \mid s \in \text{dom}(\text{Cycle}(R))) \\ \text{distance-set}(R, s) = \\ \{ d : \mathbb{N} \mid (\exists sp : \text{Span}(R)) (d = \text{distance}(sp, s)) \} \\ \text{where} \\ \text{distance}(sp, s) = \min \{ n : \mathbb{N} \mid sp^n(s) \notin \text{dom}(sp) \} \\ \text{which exists, since } sp \text{ is functional and} \\ s \in \text{dom}(loop(sp)) \end{array} \quad \text{DEF (3.15)}}$$

The set of spanning functions for a relation may be infinite; however - as the next theorem states - for a finite branching relation R the distance-sets for an element within $\text{dom}(\text{Cycle}(R))$ are finite. Hence the definition given for Rank is valid.

S

THEOREM (3.16)

Events: $\mathcal{F}(S \rightarrow S)$;
 $s : S$
 $s \in \text{dom}(\text{Cycle}(\text{UEvents})) \vdash$
 $\text{distance-set}(\text{UEvents}, s) \in \mathcal{F}(N)$

let us define -

$\text{INF} = \{ s : S \mid \text{distance-set}(\text{rel}, s) \notin \mathcal{F}(N) \}$
 $\text{rel} = \text{UEvents}$

proof : by contradiction, the assumption is

- (1) $s \in \text{INF}$
- (2) $(\forall x:S)(x \notin \text{dom}(\text{rel}) \Rightarrow \text{distance-set}(\text{rel}, x) = \{0\})$
from which we get via contrapositive
- (3) $(\forall x:S)((x \in \text{INF}) \Rightarrow x \in \text{dom}(\text{rel}))$
at least one of the finite number of states reachable from a state in INF must have an infinite distance-set, i.e. -
- (4) $(\forall x:S)(x \in \text{INF} \Rightarrow (\exists y:S)(x \text{ rel } y \wedge y \in \text{INF}))$
from (3) we get
- (5) $(\forall x:S)$
 $(x \in \text{INF} \Rightarrow (\exists y:S)(x \text{ rel } y \wedge y \in \text{INF} \wedge y \in \text{dom}(\text{rel})))$
This allows us to define a function, as follows -
 $t = (\lambda x: \text{INF})(\mu y: S)(y \in \text{INF} \wedge x \text{ rel } y)$
which gives -
- (6) $t \subseteq \text{rel} \wedge t \in S \rightarrow S \wedge t \in \text{INF} \rightarrow \text{INF}$
The function
 $\text{sp} = \tau(\text{Span}(\text{rel}) \circ t)$
is a spanning function for rel (6), and as
 $s \in \text{dom}(t)$
we get (as $t \in \text{INF} \rightarrow \text{INF}$)
 $s \notin \text{dom}(\text{loop}(\text{sp}))$
therefore
- (7) $(\exists f : \text{Span}(\text{rel})) (s \notin \text{dom}(\text{loop}(f)))$
which gives (THEOREM(3.9))
- (8) $s \notin \text{dom}(\text{Cycle}(\text{rel}))$
which contradicts the hypothesis.

3.2. A BOUNDED BUFFER.

To illustrate the use of the definitions given in the previous section we examine an *abstract design* for a small system with the following properties :

- 1) The system consumes *input* from the *environment*.
- 2) All *input* consumed will eventually appear as *output*. Furthermore, the order in which the elements appear in the *output* stream is identical to the order in which they were consumed.
- 3) Only a limited number of elements can be 'within' the system at any given time.
- 4) Necessary conditions for an *input-event* to occur are -
 - a) the event will not cause the limit to be exceeded.
 - b) the *environment* is willing to participate in such an event, i.e. there are elements to consume.
- 5) A necessary condition for an *output-event* to occur is that there are still elements, previously consumed, left 'within' the system.

This informal specification will be the basis for the *abstract design* presented in this section. In succeeding chapters formal specification techniques will be developed, but to justify these techniques we need first to examine the general properties of non-deterministic systems. It should be emphasized that the following description is not considered to be a *formal specification*. It is a *design* we propose, and we will use the tools from the previous section to investigate how this particular design relates to the informal specification given.

Given the environment -

$\text{in_env} : \text{seq}[X] ; \quad \text{Bound} : N1$

then an abstract state of the system may be described as follows :

BB _____ X _____ DEF (3.17)

$\text{unseen, buf, out} : \text{seq}[X]$

$\# \text{buf} \leq \text{Bound};$
 $\text{in_env} = \text{out} * \text{buf} * \text{unseen}$

where *unseen* describes the un-consumed input, *out* the produced output and *buf* the elements 'within' the system.

As operations on this state we propose two non-deterministically occurring events. The necessary condition for an input to occur is described in *PRE_INPUT*, the state transformations are described in *INPUT*.

PRE_INPUT $\frac{X}{BB ; \text{unseen} \neq \langle \rangle ; \# \text{buf} < \text{Bound}}$ DEF (3.18a)

INPUT $\frac{X}{\text{buf}' = \text{buf} * \langle \text{first}(\text{unseen}) \rangle ; \text{unseen}' = \text{ending}(\text{unseen}) ; \text{out}' = \text{out}}$ DEF (3.18b)

Similarly for output we have

PRE_OUTPUT $\frac{X}{BB ; \text{buf} \neq \langle \rangle}$ DEF (3.19a)

OUTPUT $\frac{X}{\text{out}' = \text{out} * \langle \text{first}(\text{buf}) \rangle ; \text{buf}' = \text{ending}(\text{buf}) ; \text{unseen}' = \text{unseen}}$ DEF (3.19b)

If we restrict ourselves to considering a particular environment, e.g. where

$$\#(\text{in_env}) = 3 \wedge \text{Bound} = 2$$

then a reasonably complete picture of the dynamic behaviour of the system can be given (FIGURE (3.20)). We define the states and the events as follows-

bb = { BB };
 events = { (λ PRE_INPUT) (μ BB')(INPUT),
 (λ PRE_OUTPUT) (μ BB')(OUTPUT) }

or using the definition (3.1)

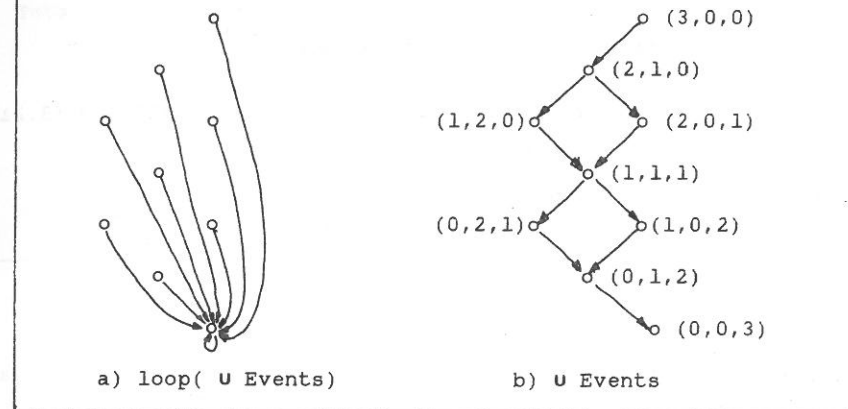
(μ ND[{{BB}}])
 (Event = { (λ PRE_INPUT) (μ BB')(INPUT),
 (λ PRE_OUTPUT) (μ BB')(OUTPUT) })

In the diagram the states are represented by triples -

(#(\text{unseen}), #(\text{buf}), #(\text{out}))

adjacent to the nodes of the behaviour-graph.

Bounded buffer Behaviour ————— FIGURE (3.20)



We now observe the following :

- 1) Independent of the starting state, it is possible to reach a terminal state in a finite number of steps, i.e.

$$(\forall b : \{BB\}) (\exists n : \mathbb{N}) (\exists b' : \{BB\}) (b \text{ Step}^n b' \wedge b' \in (\{BB\} - \text{dom}(\text{Step})))$$

where

Step = U Events

in other words the behaviour can be described as a total relation -
 $\text{dom}(\text{loop}(\text{Step})) = \{BB\}$

- 2) Independent of starting state it is guaranteed that a terminal state will be reached (the graph is loop-free), i.e.

$$(\forall \text{sp} : \text{Span}(\text{Step})) (\text{dom}(\text{loop}(\text{sp})) = \{BB\})$$

or

$$\text{dom}(\text{Cycle}(\text{Step})) = \{BB\}$$

- 3) If the system is started in a state where the buffer (buf) is empty and no output has been produced, then the system will terminate in a state where the output is the original input (in_env) - this is formulated as follows.

Given -

Effect = Cycle(UEvents)

then -

$$\text{Effect}(\{BB \mid \text{out}, \text{buf} = \langle \rangle ; \text{unseen} = \text{in_env}\}) \subseteq \{BB \mid \text{buf}, \text{unseen} = \langle \rangle ; \text{out} = \text{in_env}\}$$

The informal discussion above illustrates how the tools provided can be used to express properties about the design of non-deterministic systems.

The correctness of our design is expressed in the two following theorems-

1) The *invariant* theorem.

(the *axioms* of BB are invariant under the events)

BUF_INV _____ THEOREM (3.21)

$\vdash (\text{input} \cup \text{output}) (\{ \text{BB} \}) \subseteq \{ \text{BB} \};$
where :
 $\text{input} = (\lambda \text{PRE_INPUT}) (\mu \text{BB}')(\text{INPUT});$
 $\text{output} = (\lambda \text{PRE_OUTPUT})(\mu \text{BB}')(\text{OUTPUT})$

which is a direct consequence of LEMMA(3.22-23)

Buf_Sys_Invariant _____ LEMMA (3.22a)

BB;
 $\text{unseen} \neq \langle \rangle \wedge \#(\text{buf}) < \text{Bound};$
 $\text{buf}' = \text{buf} * \langle \text{first}(\text{unseen}) \rangle;$
 $\text{unseen}' = \text{end}(\text{unseen}); \text{out}' = \text{out}$
 $\vdash \text{BB}'$

We may use the descriptions introduced to restate LEMMA(3.22a) more compactly.

Buf_Sys_Invariant _____ LEMMA (3.22b)

PRE_INPUT; INPUT $\vdash \text{BB}'$

Buf_Sys_Invariant _____ LEMMA (3.23)

PRE_OUTPUT; OUTPUT $\vdash \text{BB}'$

2) The *total correctness* theorem

Total correctness _____ THEOREM (3.24)

Step : { BB } \leftrightarrow { BB }
 $\text{Step} = (\lambda \text{PRE_INPUT}) (\mu \text{BB}')(\text{INPUT}) \cup$
 $(\lambda \text{PRE_OUTPUT})(\mu \text{BB}')(\text{OUTPUT})$
 \vdash
 $(\text{Cycle}(\text{Step}))(\{ \text{BB} \mid \text{out}, \text{buf} = \langle \rangle; \text{unseen} = \text{in_env} \}) \subseteq$
 $\{ \text{BB} \mid \text{unseen}, \text{buf} = \langle \rangle; \text{out} = \text{in_env} \}$

proof :

Note, the following proof is conditional as the proof of (1) is assumed.
 The omitted part of the proof will be given in section 3.3.2.

let us define

Effect = Cycle(Step);

Initial = { BB | out, buf = $\langle \rangle$, unseen = in_env }

Final = { BB | unseen, buf = $\langle \rangle$; out = in_env }

- (1) $\text{dom}(\text{Effect}) = \{ \text{BB} \}$
- (2) $\text{Initial} \subseteq \text{dom}(\text{Cycle}(\text{Step}))$ (1)
- (3) $\text{Step}(\{ \text{BB} \}) \subseteq \{ \text{BB} \}$ (TH 3.21)
- (4) $\text{Effect}(\text{Initial}) \subseteq \{ \text{BB} \} - \text{dom}(\text{Step})$ (TH 3.11)
- (5) $\text{Effect}(\text{Initial}) \subseteq$
 $\{ \text{BB} \mid (\text{unseen} = \langle \rangle \vee \# \text{buf} = \text{Bound}) \wedge \text{buf} = \langle \rangle \}$
- (6) $\text{Effect}(\text{Initial}) \subseteq$
 $\{ \text{BB} \mid \text{unseen} = \langle \rangle \wedge \text{buf} = \langle \rangle \} \cup$
 $\{ \text{BB} \mid \# \text{buf} = \text{Bound} \wedge \text{buf} = \langle \rangle \}$
- (7) $\{ \text{BB} \mid \# \text{buf} = \text{Bound} \wedge \text{buf} = \langle \rangle \} = \{ \}$ (BoundeN1)
- (8) $\text{Effect}(\text{Initial}) \subseteq$
 $\{ \text{BB} \mid (\text{unseen} = \langle \rangle \wedge \text{buf} = \langle \rangle) \}$ (6), (7), (BB)
- using the invariant we get
- (9) $\text{Effect}(\text{Initial}) \subseteq \text{Final}$

Consider a slightly modified version of the bounded buffer with a *PUTBACK* facility.
 The *putback-event* is described as follows :

PRE_PUTBACK _____ X DEF (3.25a)

BB ; $\#(\text{buf}) = \text{Bound}$

PUTBACK _____ X DEF (3.25b)

$\text{unseen}' = \langle \text{last}(\text{buf}) \rangle * \text{unseen};$
 $\text{buf}' = \text{beginning}(\text{buf});$
 $\text{out}' = \text{out}$

The dynamic behaviour of the system described as follows -

BUF_with_PUTBACK $\overline{\text{X}}$ DEF (3.26)

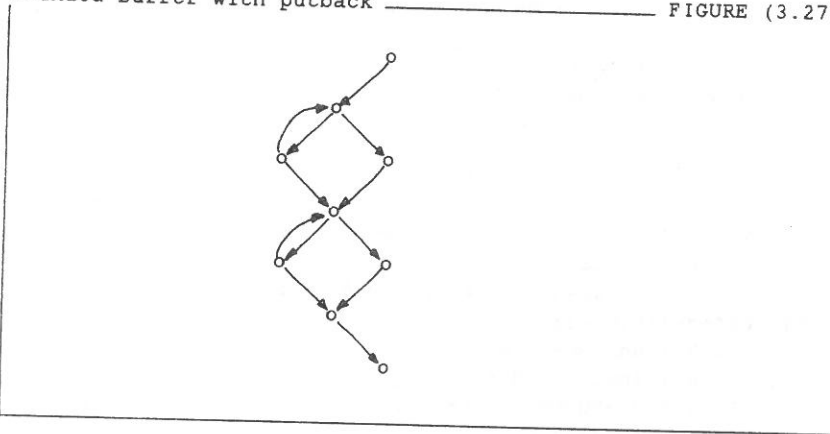
```

in_env : seq[X]
Bound  : N1
Step   : {BB} ↔ {BB}

Bound = 2 ; #( in_env ) = 3;
Step = (λ PRE_INPUT) (μ BB')(INPUT) ∪
        (λ PRE_OUTPUT) (μ BB')(OUTPUT) ∪
        (λ PRE_PUTPACK) (μ BB')(PUTBACK)
    
```

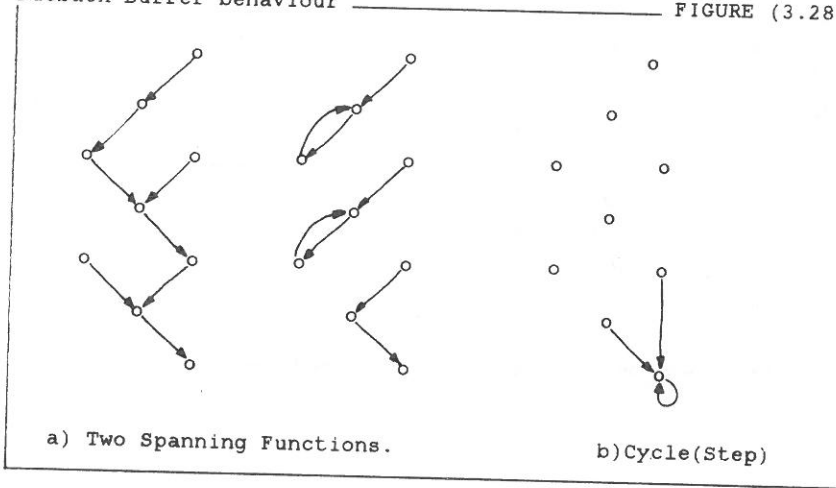
is illustrated in -

Bounded Buffer with putback $\overline{\text{X}}$ FIGURE (3.27)



This system is not guaranteed to terminate as an infinite sequence of 'input - putback' events may occur. FIGURE(3.28) illustrates two of the spanning-functions and the partial function Cycle(Step).

Putback-Buffer behaviour $\overline{\text{X}}$ FIGURE (3.28)



The dynamic behaviour of the system may be characterized formally:

$\overline{\text{X}}$ THEOREM (3.29)

```

BUF_with_PUTBACK;
f1 , f2 : {BB} → {BB};
f1 = Cycle( Step );
f2 = Loop( Step )      ⊢

dom( f2 ) = {BB};
dom( f1 ) = {BB | #(buf * unseen) < Bound}
    
```

which reads - independent of the starting state the system may possibly terminate; termination cannot be guaranteed, unless the number of elements yet to be produced is less than the bound.

3.3 FORMALISATION OF NON-DETERMINISTIC SYSTEMS.

The definitions and theorems in section 3.1 gave a formal characterization of non-deterministic behaviour.

The important theorems comprise- 1) [THEOREM(3.13): Termination] which describes the proper starting states for a system, i.e. the states from which terminal-states will be reached - 2) [THEOREM(3.11): CYCLE] which is used to further characterize the states in which a particular non-deterministic system (NDS) terminates.

It should be mentioned that termination properties are considered to be an important concept when designing systems consisting of several independent processing components. *Distributed programs* which are intended to co-operate in solving a specific problem, such as finding a root for a particular function, are expected to terminate.

Systems which ideally should run continuously - such as Operating Systems, Distributed Data Base Systems or Network Communication Systems - should be designed to terminate *only* if the system has been stopped intentionally (operator stop) or if a failure occurs within a vital component of the system, and furthermore, only when all current activity unaffected by the possible failure is properly terminated.

This approach requires that systems for which we prove termination are assumed to operate in a finite environment, i.e. the number of incoming requests (cards, OS-commands, failure-messages) is bounded.

The termination properties which are built into the design of these kinds of systems will, in the implementation, be provided by *close-down* procedures.

We use the following classical concepts to formally discuss termination properties for non-deterministic systems :

Final - states: states in which a system is intended to terminate.

Deadlock - states: terminal states, which are not foreseen Final states.

Initial - states: the states from which the system may be started.

The important categories of non-deterministic systems can be described as follows :

- 1) The unrestricted system, which may deadlock or encounter an infinite loop.

$\text{NDS} \xrightarrow{X} \text{DEF (3.30)}$ <p>Events : $\mathcal{F}(X \rightarrow X)$; Terminal : $\mathcal{F}(X)$; Initial, Final, Deadlock : $\mathcal{F}(X)$</p> <hr/> <p>Terminal = $X - \text{dom}(\cup \text{Events})$; Deadlock = Terminal - Final</p>
--

- 2) The halting NDS (include freedom from loops).

$\text{HALTING_NDS} \xrightarrow{X} \text{DEF (3.31)}$ <p>NDS[X]</p> <hr/> <p>Initial $\subseteq \text{dom}(\text{Cycle}(\cup \text{Events}))$</p>
--

intuition - THEOREM(3.13): Termination

- 3) The *totally correct* system. A system is considered to be totally correct if it can be proved to terminate (loop-free) and that all possible Terminal states, reachable from the Initial states, are acceptable Final states (deadlock-free).

$\text{WELL_HALTING_NDS} \xrightarrow{X} \text{DEF (3.32)}$ <p>HALTING_NDS[X]</p> <hr/> <p>(Cycle($\cup \text{Events}$))(Initial) \subseteq Final</p>

- 4) The very 'robust' system which behaves 'correctly' independent of the starting state, the system may compute the wrong result if it is started outside the Initial state, but it will never encounter an infinite loop or deadlock.

$\text{NON_BREAKABLE_NDS} \xrightarrow{X} \text{DEF (3.33)}$ <p>NDS</p> <hr/> <p>dom(Cycle($\cup \text{Events}$)) = X; Terminal \subseteq Final</p>
--

These properties are indeed sufficient to guarantee *total correctness* as we have :

$\text{NON_BREAKABLE_NDS} \xrightarrow{X} \text{THEOREM(3.34)}$ <p>NON_BREAKABLE_NDS \vdash WELL_HALTING_NDS</p> <hr/> <p>let Step = $\cup \text{Events}$</p> <p>proof :</p> <p>(1) dom(Cycle(Step)) = X (NON_BREAKABLE) (2) Initial \subseteq dom(Cycle(Step)) (1) (3) X - dom(Step) \subseteq Final (NON_BREAKABLE) (4) ran(Cycle(Step)) \subseteq X - dom(Step) (TH (3.11)) (5) (Cycle(Step))(Initial) \subseteq Final (3),(4) (6) WELL_HALTING_NDS (2),(5)</p>
--

3.3.1 Deadlock and the Invariant Condition.

Using the new definitions the statement of correctness [THEOREM(3.24): Bufsys] for the design of the small Buffer System in section 3.2 may now be reformulated (NB. THEOREM(3.35) is stronger than THEOREM(3.24)).

$\text{BUFSYS} \xrightarrow{X} \text{THEOREM(3.35)}$ <p>in_env : seq[X] ; Bound : N1 NDS[{BB[X]}]; Events = { ($\lambda \text{PRE_INPUT}$)($\mu \text{BB}'$)(INPUT), ($\lambda \text{PRE_OUTPUT}$)($\mu \text{BB}'$)(OUTPUT) } Initial = { BB out=$\langle \rangle$; buf=$\langle \rangle$; unseen=in_env } Final = { BB buf=$\langle \rangle$; unseen=$\langle \rangle$; out=in_env } \vdash NON_BREAKABLE_NDS[{BB}]</p>
--

The proof of THEOREM(3.24) illustrates the importance of the *Invariant theorem* in proving WELL_HALTING properties.

If the small buffer system is started outside Initial-states (e.g. with a non-empty buffer) then the Final-states are still the only Terminal-states.

however the function computed by the system is no longer the intended identity for sequences. This property, which is the *NON_BREAKABLE* property, is directly related to the *invariants*, which constrain the states.

If the abstract state (DEF(3.17) - BB) had been -

$\text{BB2} \xrightarrow{\quad X \quad} \text{DEF (3.36)}$ $\text{unseen, buf, out : seq}[X]$ $\#(\text{buf}) \leq \text{Bound}$
--

i.e. without the important invariant

$$\text{in_env} = \text{out} * \text{buf} * \text{unseen}$$

which says that activities performed by the system preserve the ordering of the data within the data-stream and do not introduce extraneous elements, then the corresponding NDS[BB] (constructed as in THEOREM(3.35) Bufsys) could have been proved to be a *WELL_HALTING_NDS* by the invariant theorem and observing that the invariant holds in the initial state.

However, such a system does not have the *NON_BREAKABLE_NDS* property, as a *Final* state cannot be reached if the system is started in a state where

$$\text{in_env} \neq \text{out} * \text{buf} * \text{unseen}$$

In general we should always construct *NON_BREAKABLE_NDS* designs. If a design can be proved to have the *WELL_HALTING* property only, then there may exist states - satisfying the invariant - from which the system may deadlock or loop indefinitely. This is indeed undesirable. Later modifications to the system, (e.g. addition of operations which can be proved to keep the invariant, but may take the system to new states not previously reachable from *Initial*) may affect the behaviour of remaining 'unmodified' parts of the system - and can cause deadlocks.

A 'breakable-nds' design shows that there are properties of the system which the designer has not observed or considered.

There are several techniques for removing deadlock-possibilities.

Given -

Final : $\mathcal{P}(\{S\})$
 Step : $\{S\} \leftrightarrow \{S\}$
where
 $\{S\} - \text{dom}(\text{Step}) \not\subseteq \text{Final}$

then there is a possibility for deadlock. This possibility must be detected and removed by deciding whether to -

- 1) accept the deadlock, i.e. extend the *Final* - states.
- 2) provide a deadlock detection solution, i.e. extend the relation *Step* by adding one or more new events.
- 3) remove the deadlock possibility by limiting the abstract space, i.e. add invariants.

as both -

$$\{S\} - \text{dom}(\text{Step} \cup \text{new-function}) \subseteq \{S\} - \text{dom}(\text{Step})$$

and

$$\{S | \text{new-invariant}\} - \text{dom}(\text{Step}) \subseteq \{S\} - \text{dom}(\text{Step})$$

the changed design of the system may well be provable to be *NON_BREAKABLE*.

3.3.2. Termination and the Variant Function.

The properties of a *HALTING_NDS* only guarantee termination provided the starting state was within the predefined *Initial*-states. As mentioned earlier we should strive to design *NON_BREAKABLE_NDS*, i.e. systems which always terminate.

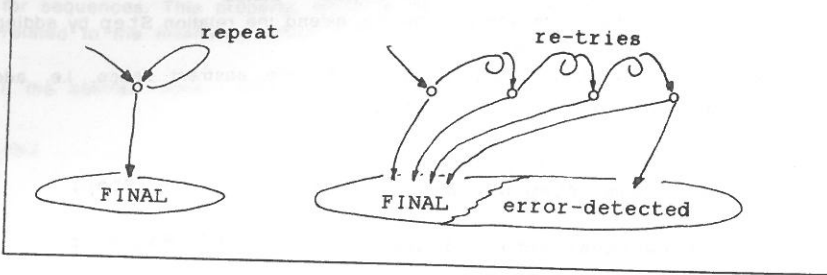
There is an important group of system components where a 'repeat' mechanism is required; e.g. a communication system should include facilities for *re-transmission* of messages which, due to unreliable communication-lines, have been lost; and a disc-handler should provide a *re-try* mechanism which will be used after the occurrence of a disc-read check-sum error.

However this 'repeat' facility should not prevent the designer from constructing *totally correct* and non-breakable systems.

Consider the 'behaviours' illustrated in FIGURE(3.37).

- a) A partially correct system, which 'only' fails to terminate WHEN the unexpected happens (loops infinitely).
- b) The 'repeat' mechanism from system a) is replaced by a well-known mechanism, in which only a limited number of re-tries is performed before the system is stopped within the new, extended *Final* - state. This mechanism will behave identically for transient errors, but persistent errors will be detected. Note, this modification is done using a similar techniques as those applied in order to prevent deadlock (see previous page). This time in order to satisfy a termination-axiom-

$$\text{dom}(\text{Cycle}(\text{a-relation})) = \text{all-states}$$

Repeat Mechanisms FIGURE(3.37)

The loop-free property of a NON_BREAKABLE_NDS, requires a proof of-

$$\text{dom}(\text{Cycle}(\text{Step})) = X$$

where $\text{Step} : X \leftrightarrow X$

is a total relation.

i.e. that the relation $\text{Cycle}(\text{Step})$

The following theorem gives a proof method for finite branching relations-

VARIANT X THEOREM(3.38)

$$\begin{aligned} \text{NDS} \quad & \vdash \\ \text{dom}(\text{Cycle}(\text{Step})) &= X \iff \\ (\exists v : X \rightarrow \mathbf{N})(\forall x, x' : X) & (x \text{ Step } x' \Rightarrow v(x') < v(x)) \\ \text{where} & \\ \text{Step} &= \cup \text{Events} \end{aligned}$$

proof :

- a) If $\text{Cycle}(\text{Step})$ is total then a variant function exists :
- (1) $\text{dom}(\text{Rank}) = X$ (DEF(3.14))
 - (2) $(\forall x, x' : X | x \text{ Step } x') (\text{Rank}(x') < \text{Rank}(x))$ (DEF(3.15))
i.e. RANK is a variant function.
- b) If a variant function for the relation exists then this variant function is indeed a variant function for all spanning-functions for Step, i.e.
- (1) $(\forall \text{sp} : \text{Span}(\text{Step})) (\forall x : \text{dom}(\text{sp})) (v(\text{sp}(x)) < v(x))$
 - (2) $(\forall x : X) (\forall \text{sp} : \text{Span}(\text{Step})) (x \in \text{dom}(\text{loop}(\text{sp})))$ (next LEMMA)
 - (3) $\text{dom}(\text{Cycle}(\text{Step})) = X$ (TH (3.9))

The well-known theorem on the importance of the variant function for iterations of partial functions -

VARIANT X LEMMA (3.39)

$$\begin{aligned} f : X &\rightarrow X \\ v : X &\rightarrow \mathbf{N} \\ (\forall x : \text{dom}(f)) & (v(f(x)) < v(x)) \\ \vdash & \\ (\forall x : X) & (\exists n : \mathbf{N}) (\exists y : X) (y = f^n(x) \wedge y \notin \text{dom}(f)) \end{aligned}$$

proof :

(1) for an arbitrary $x : X$ we define the following partial function:

$$\begin{aligned} g_x &= (\lambda n : \mathbf{N} | f^n(x) \in \mathbf{N}) (v(f^n(x))) \\ &\text{because } (0, v(x)) \in g_x \text{ we can define :} \\ m_x &= \min(\text{ran } g_x) \\ n_x &= \min g_x^{-1}(\{m_x\}) \\ &\text{hence, we can define} \end{aligned}$$

(2) $y_x = f^{n_x}(x)$

under the assumption that

$$y_x \in \text{dom}(f)$$

we get (as v is a decreasing function)

$$g_x(n_x + 1) < \min(\text{ran } g_x),$$

which is a contradiction, therefore

(3) $y_x \notin \text{dom}(f)$

hence according to (1), (2) and (3) we get

$$(\forall x : X) (\exists n : \mathbf{N}) (\exists y : X) (y = f^n(x) \wedge y \notin \text{dom}(f))$$

Taking the proof-method for proving freedom from loops into account, we can now reformulate our definition of NON_BREAKABLE_NDS -

NON_BREAKABLE_NDS X DEF (3.40)

$$\begin{aligned} \text{Events} & : \mathcal{F}(X \rightarrow X) \\ \text{Initial, Final} & : \mathcal{P}(X) \\ \text{Var} & : X \rightarrow \mathbf{N} \end{aligned}$$

$$\begin{aligned} (X - \text{dom}(\cup \text{Events})) &\subseteq \text{Final} \\ (\forall f : \text{Events})(\forall x : \text{dom}(f)) & (\text{Var}(f(x)) < \text{Var}(x)) \end{aligned}$$

Example:

The statement of correctness for the small Buffer System (section 3.2) may now be formulated-

Buf_Sys $\frac{\quad}{X}$ THEOREM (3.41)

```

in_env  : seq[X]
Bound   : N1
Events  : F( {BB} → {BB} )
Initial,
Final   : P( {BB} )
Var     : {BB} → N
Events  = { (λ PRE_INPUT)(μ BB')(INPUT),
            (λ PRE_OUTPUT)(μ BB')(OUTPUT) }
Initial = { BB | out,buf = <> }
Final   = { BB | unseen,buf = <> }
Var     = (λ BB) ( 2 * #(unseen) + #(buf) )
├
NON_BREAKABLE_NDS[ {BB} ]

```

Note! the proof of the total correctness of the Buffer System (which was conditional in THEOREM(3.24)) may now be completed by proving that the given variant function is indeed a variant function for all participating events.

We must prove-

Buf_Sys_Variant $\frac{\quad}{X}$ LEMMA (3.42)

```

PRE_INPUT; INPUT  ┆
(2*#(unseen')) + #(buf') < (2*#(unseen)) + #(buf);

```

Buf_Sys_Variant $\frac{\quad}{X}$ LEMMA (3.43)

```

PRE_OUTPUT; OUTPUT ┆
(2*#(unseen')) + #(buf') < (2*#(unseen)) + #(buf);

```

LEMMA(3.22) and LEMMA(3.23) guarantee that

$$\text{Events} \in F(\{BB\} \rightarrow \{BB\})$$

We can prove the deadlock-free property -

$$(X - \text{dom}(\cup \text{Events})) \subseteq \text{Final}$$

as in THEOREM(3.24) or by proving -

Buf_sys_termination $\frac{\quad}{X}$ LEMMA (3.44)

```

in_env  : seq[X];
Bound   : N1;
BB;
(unseen=<> ∨ #(buf)=Bound) ∧ buf=<>;
┆
( unseen=<> ∧ buf=<> ∧ out=in_env )

```

which expresses that the Buffer System terminates within the acceptable final states only.

The method illustrated here for proving non-deterministic designs correct can be summarised as follows:

- 1) State the NON_BREAKABLE properties as in THEOREM(3.41).
- 2) Prove the Invariant Lemmas. (see LEMMA(3.22-23)).
- 3) Prove the Variant Lemmas. (LEMMA(3.42-43)).
- 4) Prove the Deadlock-Lemma. (LEMMA(3.44))

Even for small systems it may be difficult to find a Variant Function. The reason may be that the *Invariants* which constrain the state-space are insufficient, whence we are faced with another difficult task - to find proper Invariants.

However, it is easy to prove or disprove that a particular function is a Variant Function for a given set of partial functions.

4. SPECIFICATIONS OF DISTRIBUTED SYSTEMS.

The remaining parts of this thesis present a new approach for specifying *DISTRIBUTED SYSTEMS*.

Distributed systems are systems which are composed of several concurrently operating processing components. As mentioned in chapter 3 these systems fall into the category of *non-deterministic systems*. The formal framework which we developed in the previous chapter will be applied to the analysis of distributed systems, and we will see (in section 4.4) how the difficulties in proving termination (i.e. finding a variant function) are reduced owing to the method used for specification of the behaviour of such systems.

4.1. MOTIVATION FOR DISTRIBUTION.

A specification of a distributed system must document the decisions taken with respect to the distribution of data- and control-structures onto independent processing units. The reasons for designing distributed systems may vary:

- a) The decisions may be directly dictated by demands with respect to *security* (e.g. certain data must only be available at a specific geographical location) or *reliability* (e.g. the introduction of redundancy).
- b) The decisions may also be made by the designer in order to satisfy *efficiency requirements* (e.g. 1) in a computer network the response time for accessing data is dependent on the 'distance' to the location of the accessed data, and 2) the throughput of a system can be increased by utilising a *pipelined* processor arrangement.)
- c) *Maintenance* and *modification* of running systems are other factors which may affect decisions with regard to the physical distribution of control- and data-structures.
- d) Finally, a decision to provide a distributed system for solving a specific task may simply be made because of the availability of 'off the shelf' hardware- or software-units (e.g. special purpose microprocessors).

The approach presented here for documenting the distribution is based on the classical idea [Parnas,20] of specifying a system in terms of descriptions of its subcomponents and the relationship between those subcomponents. In our approach each component of the *description* corresponds to a single component or a set of components in the *realisation*.

4.2. THE INDEPENDENT MODULE.

When describing a subcomponent of a system it is essential to draw a distinction between its externally observable behaviour and its internal composition. If we have a complete description of a component's *external connections* and a description of the behaviour along those connections there is nothing more of interest we can learn by looking at the inside of the component. The main problem, however, is to describe the external behaviour entirely in terms of those connections, i.e. without any reference to the internal structure. For example, the specification of external behaviour of the *Buffer* from section 3.2 does not qualify as it involves an abstract structure (*buf*) which represents the inside of the described component.

4.2.1. The Classical Independent Module.

The classical independent module can be described through an *input output* relationship, i.e. in terms of its external behaviour only. A specification of a module through an input-output relationship gives a complete description of the behaviour of a realisation of such a module because we assume that any realisation follows two principles:

- 1) there are only two kind of external connections - *input* and *output* - and these connections are *activated* alternately.
- 2) any output value produced is exclusively derivable from the module's functional definition and from the value of the most recent input presented, i.e. the module is *stateless*.

4.2.2. The Process Module.

The modules under consideration here do not follow either of these simple principles. For the modules in question the ability to participate in a communication along an external connection is not governed by a simple data-independent control mechanism, and a value communicated along a connection is not determined solely by the most recent values communicated along another set of connections. At any point in time both the ability to communicate and the actual value communicated along a particular connection can be dependent on information in the *past history* of the behaviour along any set of external connections. The distinction between the classical independent module and the modules we are interested in - hereafter called *Process Modules* - can be compared to the distinction between combinatorial logic circuits (e.g. AND-gates) whose output depends only on the most recent input, and sequential logic circuits whose output depends on past input (the history of the input activity), e.g. the output of a flip-flop is determined by the parity of the number of input-signals previously transmitted to it.

NB, a reasonable realisation of a module whose behaviour is dependent on its past behaviour memorises only a minimal set of facts about the past only - e.g. the flip-flop has a 1-bit memory which represents the *parity of the number of input signals*.

The independent module's behaviour can be specified by giving:

- 1) The type of input and output parameters.
- 2) A functional definition of the output produced, i.e.

$$\text{output} = F(\text{input})$$

The process module's behaviour can be specified by giving:

- 1) The type of the objects which can be communicated along its external connections.
- 2) A set of axioms involving one or more of the past histories of the behaviour along those external connections.

Example:

We will give a specification of a process (P1) with two external connections (*in* and *out*) which for each element (*x*) communicated along the input connection (*in*) will produce a single result (*F(x)*) along the output connection (*out*) before any new input can be consumed.

Note the similarities with the assumptions made for a realisation of independent memoryless modules.

Given $F : X \rightarrow Y$, we will define the behaviour of our process by, 1) modelling the history of the communications along the input connection as a *sequence* over elements of type *X*, 2) modelling the history of the communications along the output connection as a *sequence* over elements of type *Y* and 3) giving two axioms involving these abstractions. We have-

P1	X Y	DEF(4.1)
$\text{in} : \text{seq}[X];$ $\text{out} : \text{seq}[Y]$		
$\# \text{in} - \# \text{out} \leq 1;$ $\text{out} \subseteq F \circ \text{in}$		

where the first axiom restricts the number of computations P1 can carry out 'simultaneously'. In this example only one computation can be carried out at a time, hence the input connection will be *blocked* after a value has been passed along the input connection and the corresponding result not yet produced along the output connection. The second axiom states that all values produced by P1 can be obtained by applying *F* to a previously received input element.

Assume that we have observed

a, b, c, and d

being communicated along the in connection in that order, then

$$\text{in} = \langle a \ b \ c \ d \rangle$$

The history of the communication along the out connection must then be

$$\langle F(a) \ F(b) \ F(c) \rangle$$

or

$$\langle F(a) \ F(b) \ F(c) \ F(d) \rangle$$

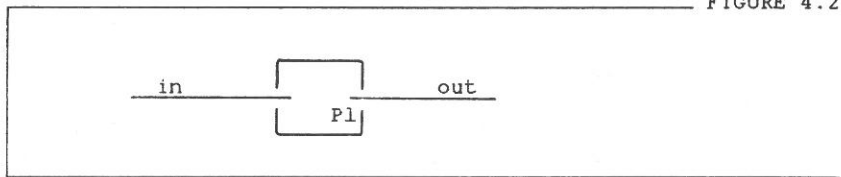
If

$$\text{out} = \langle F(a) \ F(b) \ F(c) \rangle$$

further input must be rejected as the occurrence of an input will invalidate the first of the axioms. However P1 may produce F(d), as this event will not invalidate any of the given axioms.

We shall in the following use a graphical notation to give an informal overview of the interaction between the process-modules of a system.

The process P1 will be represented by



which may be interpreted either as a box representing process component which realises module P1, or simply as an illustration of the fact that P1 is a description which involves axioms on in and out only.

4.2.3. The Readiness Principles.

A realisation of a process module must - just like a sequential realisation of a classical independent module - follow certain principles. As a process module contains information concerning the scheduling of the activities along its external connections the principles take a different form from those for the independent module.

The principles are:

- 1) At any point in time the realisation of a process-module must not participate in a communication (i.e. extend a history) if the occurrence of this communication would invalidate the axioms given (for the histories of the activities along the external connections). We might say - 'only communications which are permitted to happen can happen'.

- 2) If an occurrence of a communication along a specific connection can take place without invalidating the external invariants then the realisation of the module must eventually, and without further activity along other connections, be willing to participate in this particular communication. We might say - 'an event which is expected to occur can eventually take place.'

No attempts will be made in this thesis to formalise these principles, for example by giving proof-rules for a specific realisation technique. We will concentrate on discussing the correctness, or rather the consistency, of specifications. However in order to illustrate how these principles can be used in statements about the correctness of an implementation we will give a simple realisation of module P1.

A realisation is outlined using the tasking facilities of the programming language *Ada*.

P1 X Y EXAMPLE (4.3)

```

task body P1 is
  x' : X ; y' : Y ;
  function F(x : in X) return Y is ... ;
begin
  loop
    accept input(x : in X) do x':=x ; end ;
    y' := F(x') ;
    accept output(y : out Y) do y:=y' ; end ;
  end loop ;
end ;

```

This is a realisation of the process-module P1 because

- 1) before the occurrence of an *input* we can assert that

$$\#in - \#out = 0$$
 hence an occurrence of an input will not invalidate the axioms of the specification P1.
- 2) before the occurrence of an *output* we can assert that

$$\#in - \#out = 1$$
 ,

$$y' = F(\text{last}(\text{in})) \quad \text{and}$$

$$\text{out} = \text{beginning}(\text{in})$$
 hence the occurrence of the output will not invalidate the axioms of P1 as the occurrence of output will update the history of communications along the out connection as follows :

$$\text{out} := \text{out} * \langle F(\text{last}(\text{in})) \rangle$$
- 3) after the occurrence of an input we have

$$\#in - \#out = 1$$
 hence we expect only an output to happen, which in fact the realisation will eventually be willing to perform.

- 4) after initiation of the task, as well as after the occurrence of an output, an input is expected to happen - which the realisation will eventually be willing to participate in.

1) and 2) ensure that the first of the principles is followed, while 3) and 4) ensure that the second principle is followed.

The example presented here as an introduction to the approach of specifying distributed systems could have been described simply by giving its functional specification, and by assuming that the module would be realised as a conventional *memoryless* software module. However, as we shall see in the next section, a small modification will give a specification of a module with a different dynamic behaviour - a behaviour which cannot be described through a simple input output relationship.

4.3. DISTRIBUTED SPECIFICATIONS AND DISTRIBUTED REALISATIONS.

Consider the specification -

P2	X Y	DEF (4.4)
		in : seq[X]; out : seq[Y]
		$\delta \leq 2$; out $\subseteq F \circ$ in
		where
		$\delta = \#in - \#out$

which is a specification of a process which computes the function F. However in contrast to the process P1 two computations can be carried out simultaneously. A realisation of this module does not follow the principles given for the classical independent module, e.g. whenever

$$\#in - \#out = 0$$

the process may participate in two consecutive inputs.

A realisation could be -

P2	X Y	EXAMPLE (4.5)
		<pre> task body P2 is i, o, δ : INTEGER := 0; queue : array (0..1) of X; function F(x : in X) return Y is ...; begin loop select when δ < 2 => accept input(x : in X) do queue(i) := x; end; δ := δ+1; i := (i+1) mod 2; or when δ ≠ 0 => accept output(y : out Y) do y := F(queue(o)); end; δ := δ-1; o := (o+1) mod 2; end select; end loop; end P2; </pre>

or if we can decompose the computation of F - i.e. we have

$$F = H \circ G$$

we might decide to provide the following *distributed* realisation -

P2	X Z Y	EXAMPLE (4.6)
		<pre> task body P2_1 is t : X; z : Z; function G(x : in X) returns Z is ...; begin loop accept input(x : in X) do t:=x ; end; z := G(t); P2_2.internal_com(z); end loop; end P2_1; task body P2_2 is t : Z; y' : Y; function H(z : in Z) return Y is ...; begin loop accept internal_com(z : in Z) do t:=z ; end; y' := H(t); accept output(y : out Y) do y:=y' ; end; end loop; end P2_2; </pre>

EXAMPLE(4.5) gives a *sequential non-distributed* realisation of DEF(4.4). The external operations (input and output) will always happen in sequence although the actual order in which they happen is unpredictable (the run-time environment

determines the order). EXAMPLE(4.6) gives a *non-sequential distributed* realisation of the same specification. The actual occurrence of input and output cannot be related (in time) unless there are some (unknown) constraints which synchronise the part of the environment which produces input with the part which consumes output.

Consider the *distributed* realisation of P2 - In the development (construction) of the task P2_1 it is of no importance how the 'inside' of P2_2 is constructed. The only part of P2_2 which is of interest in developing P2_1 is the common interface between these two processes. P2_1 can be viewed as a realisation of the specification-

P1_1	X Z	DEF (4.7)
$in : seq[X]; \quad com : seq[Z]$		
$\#in - \#out \leq 1; \quad com \subseteq G \bullet in$		

and P2_2 can be viewed as a realisation of the specification-

P2_2	Z Y	DEF (4.8)
$com : seq[Z]; \quad out : [Y]$		
$\#com - \#out \leq 1; \quad out \subseteq H \bullet com$		

The following *compound* specification

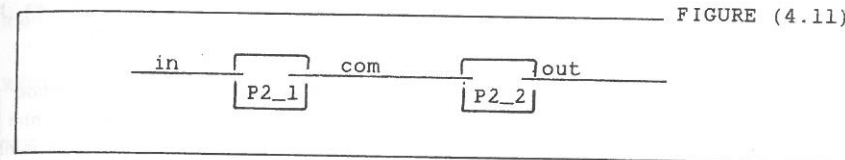
P2-Dis	X Z Y	DEF (4.9)
$P2_1[X,Z] ; P2_2[Z,Y]$		

which expands into

P2-Dis	X Z Y	DEF (4.10)
$in : seq[X]; \quad com : seq[Z]; \quad out : seq[Y]$		
$\#in - \#com \leq 1; \quad com \subseteq G \bullet in;$		
$\#com - \#out \leq 1; \quad out \subseteq H \bullet com;$		

is a *distributed specification* as the set of axioms of the description can be divided into two groups, a group which involves the in and com connections only and a group which involves the com and out connections only.

We can use our graphical notation to illustrate the distribution-



This illustration has an obvious alternative interpretation: two processes communicating along the connection com.

We may attempt to prove that the *distributed realisation* (EXAMPLE(4.6)) is a valid implementation of its *non-distributed specification*, - which means that we must demonstrate that the realisation follows the *readiness principles* with respect to its specification DEF(4.4). Such a proof will not be straightforward since the correctness of, for example, the task P2_1 (i.e. a proof that it behaves correctly along its connection in) depends on the correctness of the task P2_2. Furthermore even if we assume the correctness of P2_2 a proof of P2_1 cannot be carried out without referring to the internal behaviour of the realisation of P2_2. This is indeed undesirable as 1) independent development of subcomponents is impossible and 2) the future modification of the local mechanisms of any module may invalidate global proofs or proofs for other modules.

Proofs should be localised, i.e. the tasks P2_1 and P2_2 (from EXAMPLE(4.6)) should be verified against the specifications DEF(4.7) and DEF(4.8) respectively. However proving that the two tasks are 'locally' correct does not guarantee that the tasks - when running in parallel - realise the non-distributed specification DEF(4.4).

Such a guarantee can be gained by proving that the distributed specification DEF(4.9) is a *valid decomposition* of the non-distributed specification DEF(4.4). A decomposition is valid if, from the sole assumption that each process-component follows the *readiness principles* with respect to its local specification, it can be proved that the processes when co-operating will follow the *readiness principles* with respect to the global specification. It is the designer's responsibility to provide such a proof.

Decomposition (or distribution) of specifications will be treated in detail in section 6 of this thesis.

Obviously the correctness of the behaviour of a system which consists of interacting processes relies on more than just the proofs made with respect to decomposition at the abstract level. We must have a guarantee that the processes of a realisation are properly interfaced, i.e. co-operating processes must agree on a common communication protocol. Furthermore we must be certain that the interprocess control mechanisms provide an appropriate service. The communication protocol selected by the *implementor* and the mechanisms for communication provided by the underlying machinery must guarantee that, a) a communication can only take place between processes if all involved processes are willing to participate, and b) at least one of the 'currently' possible communications eventually takes

place. These aspects of the correctness of a system reflect the obvious assumption a designer makes about future realisations; they will be referred to as the *connectivity assumptions* and will not be treated further in this thesis.

The techniques presented in this section provide the designer with a method for describing dynamic behaviour at an abstract level, i.e. it is possible for him to express WHAT (i.e. which operations) will take place WHEN without determining the HOW. Furthermore by *distributing* his specification (*partitioning* his design) he can document his decisions about WHERE the operations will be performed. In section 4.5 we will demonstrate how such *distributed specifications* can be used as a basis for the analysis of dynamic behaviour.

4.4. ENVIRONMENTS FOR SYSTEMS.

A precise description of the environment in which a system is expected to operate is as important as a description of the subcomponents of which the system consists. The environment as well as any subcomponent of a system must obey certain well-documented rules. The behaviour of the environment for a system can be described as a separate independent system, i.e. as a set of assertions involving the past histories of the communications along the connections to the system for which it serves as an environment. The environment may consist of several independent parts and they should, if this is the case, be described using a distributed formulation.

A description of a system and an environment can be 'joined' together - thereby giving a complete description of the activity along the connections between the system and its environment.

As mentioned in section 3.3. we will investigate the behaviour of systems in *finite environments* only, i.e. we assume that the part of the environment which initiates new activities within the system only performs a finite number of requests. Such an 'inputting' environment must cease to communicate after having participated in a finite number of communications.

In the following we will give examples of environments, some of which will be used as environments for the systems presented in chapter 5.

4.4.1. A Requesting Environment.

Consider an environment which issues commands along a single connection to a system. The environment never commits itself to one particular command, but will be capable of 'participating' in any command which satisfies certain criteria. The environment promises to issue a limited number of requests.

Given a set of commands R

and a limit $L : \mathbb{N}$

then we can formalise such an environment as follows

REQUESTER	R	DEF (4.11)
requests : seq[R]		
#requests $\leq L$		

Example:

REQUESTER may describe a human user of an interactive operating system who requests certain services - save-a-file, edit-a-file, send-a-message, etc. The user communicates with the *front-end* of the command interpreter, which ensures that only intelligible commands are passed between the user and the operating system.

Given intelligible_cmds : $\mathcal{P}(R)$, then the description of the front-end module might be

FRONT END	R	DEF (4.13)
requests : seq[R]		
ran(requests) \subseteq intelligible_cmds		

While the FRONT-END ensures that unintelligible commands are rejected, the user (REQUESTER) promises that his session is finite, i.e. when
#requests = L
he will not issue any more requests.

4.4.2. A Selfcommitting Environment.

Consider a component of an environment which communicates a sequence of data-elements or requests along a single connection. At any point in time this component commits itself to participate in one particular communication only.

Given

in_env : seq[E]

then the environment which transmits elements from ran(in_env) along a connection in in the order they appear in the sequence in_env, can be described as-

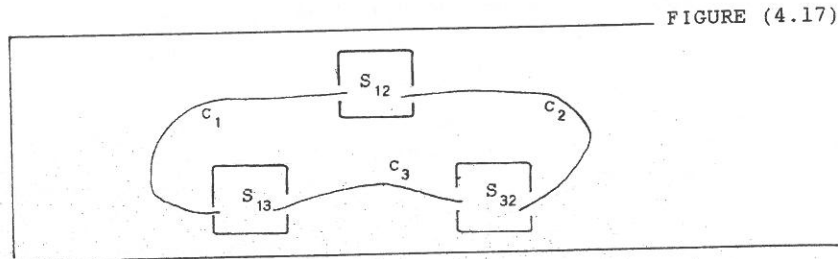
PRODUCER	E	DEF (4.14)
in : seq[E]		
in \subseteq in_env		

Operations within such a system may happen in *parallel*, simply because operations are unrelated.

In order to discuss the dynamic properties of distributed systems we base our dynamic model on the *OBSERVATIONS* of the behaviour. Without any loss of generality observations can be sequentialised; hence we can apply the formal framework developed in chapter 3.

Example:

Consider a description of a system with 3 visible connections-



The information which can be externally observed to change is-

CHAN	M	DEF (4.18)
$c_1, c_2, c_3 : \text{seq}[M]$		

where M denotes the messages communicated between the modules, and c_1 , c_2 and c_3 denote the history of observations along the visible connections.

S_{ij} is the specification of a module which communicates along c_i and c_j . We have a guarantee (if the realisation is correct) that the communications along c_i and c_j will not invalidate S_{ij} . Hence if the system is started in a state satisfying S_{12} , S_{32} and S_{31} then these will be satisfied in any state the system can ever reach. The system's observable states can therefore be described as-

SYSTEM	M	DEF (4.19)
$\text{CHAN}[M]$		
$S_{12}; S_{23}; S_{31}$		

When a communication takes place, a *single message will be appended to an existing history*, hence an occurrence of a communication along the visible connection c_1 can be described as-

M	DEF (4.20)
$\text{hist}_1 : M \rightarrow (\{ \text{CHAN} \} \rightarrow \{ \text{CHAN} \})$	
$(\forall m:M)$ $\text{hist}_1(m) =$ $(\lambda \text{CHAN})(\mu \text{CHAN}') (c_1' = c_1 * \langle m \rangle; c_2' = c_2; c_3' = c_3)$	

Furthermore, if the system is started in a state satisfying the axioms of *SYSTEM* then only communications keeping the axioms invariant will take place (according to the *readiness principles* in section 4.3.2). Hence we can refine our description of the communication along connection c_1 by *restricting its domain and co-domain* as follows-

M	DEF (4.21)
$\text{obs}_1 : M \rightarrow (\{ \text{SYSTEM} \} \rightarrow \{ \text{SYSTEM} \})$	
$(\forall m:M) (\text{obs}_1(m) = (\text{hist}_1(m) \uparrow \{ \text{SYSTEM} \}) \downarrow \{ \text{SYSTEM} \})$	

Assuming that the number of messages is finite i.e. -

MES : $F(M)$

then the total behaviour of *SYSTEM* can now be described as-

OBSERVATIONS	M	DEF (4.22)
$\text{Events} : F(\{ \text{SYSTEM} \} \rightarrow \{ \text{SYSTEM} \})$		
$\text{Events} = \cup \{ \text{obs}_1(m) \mid m:MES \} \cup$ $\cup \{ \text{obs}_2(m) \mid m:MES \} \cup$ $\cup \{ \text{obs}_3(m) \mid m:MES \}$		

So far we have a description of a state (DEF(4.19)) and a description of the transformations of that state (DEF(4.22)). By giving a description of the acceptable final states we can discuss the *NON-BREAKABLE* properties as stated in DEF(3.33) or DEF(3.40).

Let

Step : $\{ \text{SYSTEM} \} \leftrightarrow \{ \text{SYSTEM} \}$

Step = \cup Events

and

Final : $F(\{ \text{SYSTEM} \})$;

Initial : $F(\{ \text{SYSTEM} \})$

then *termination* is guaranteed if we can find a variant function, i.e. if

Var : {SYSTEM} → N
 $(\forall s, s' : \{SYSTEM\}) (s \text{ Step } s' \Rightarrow \text{Var}(s') < \text{Var}(s))$

and the system will be *deadlock-free* if

$\{SYSTEM\} - \text{dom}(\text{Step}) \subseteq \text{Final}$

NB.

- 1) Both the initial and the final state of the system are expressed in terms of the history of the visible behaviour.
- 2) The initial state for most of these systems will be a state where no communication has taken place, i.e. all history components of the state are *empty sequences*. For SYSTEM we have-

Initial = {SYSTEM | $c_1, c_2, c_3 = \langle \rangle$ }

- 3) The proof of termination relies on the existence of a variant function; for dynamic behaviour based entirely on the history of observations it is normally straightforward to find a variant function.
 E.g. for the system presented here we have-

inc_func = $(\lambda \text{SYSTEM})(\#c_1 + \#c_2 + \#c_3)$

which obviously satisfies

$(\forall s, s' : \{SYSTEM\})$
 $(s \text{ Step } s' \Rightarrow \text{inc_func}(s') > \text{inc_func}(s))$

hence the existence of a decreasing variant function relies on the existence of an *upper limit* for the total number of communications which can take place within the system. The existence of the upper limit will normally rely on an *assumption* about the environment.

- 4) A discussion of the deadlock properties of a system relies on the *consistency* of the decisions taken with respect to the dynamic behaviour (i.e. the *predicates* of the *histories*) and the definition of the final state. The next chapter will demonstrate how the *consistency check* is carried out.

5. EXAMPLES OF SPECIFICATIONS.

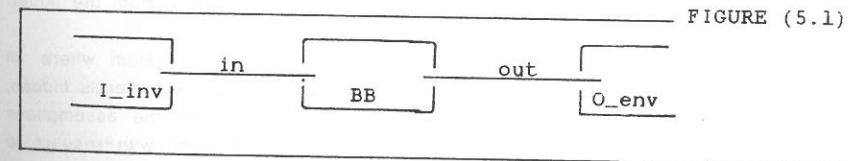
The approach presented in the previous chapter is intended to be applicable to real-life engineering problems in the area of *distributed* computing. It is not, however, intended to be suitable for examining all aspects of big distributed systems. The properties which will be documented and analysed are only those which concern co-operation between independent process components, that is, the approach will be used to give a precise formal description of the decisions made with respect to the distribution of data- and control-structures between the process components on the one hand, and a precise formal description of the interaction between these components on the other.

In this chapter the approach is illustrated through a set of simple examples. In order to compare the new approach, where behaviour is described through predicates on the histories of observed communications, with the more traditional approach, where behaviour is described through a relation (or a set of partial functions) over an abstract state, section 1 gives a new specification of the Bounded Buffer (from section 3.2). The succeeding three sections will discuss miscellaneous systems (buffers and queues) and concepts (mutual exclusion, deadlock etc.) which are of concern when designing distributed systems. The final section will give a specification of a communication network.

5.1. A BOUNDED BUFFER.

The informal description of the bounded buffer given in the beginning of section 3.2. was used to justify the design of a small system which was described in terms of 1) an *abstract state* (BB:DEF(3.17)) representing the inside of a Bounded Buffer through a *buf* component and 2) two *state transformation* functions (DEF(3.18-19)). Using the new approach we can give a description which replaces the informal specification in section 3.2.. We can formally describe a bounded buffer without suggesting an internal state for the buffer and without (misleadingly) assuming that only two *indivisible* operations make up the *dynamics* of the system.

The Bounded Buffer system can be illustrated as follows-



i.e. the system consists of three independent operating components, two of which belong to the *environment*.

We have-

- 1) The input environment is described as-

$$I_env[X] \equiv [in:seq[X] \mid \#in \leq Max] \quad DEF(5.2)$$

for some
 $Max \in \mathbb{N}1$

We assume that the inputting environment produces only a finite number of elements. (we might have used REQUESTERDEF(4.11)).

- 2) The buffer-system is described as-

$$BB[X] \equiv [in,out:seq[X] \mid out \leq in; \#in - \#out \leq B] \quad DEF(5.3)$$

for some
 $B \in \mathbb{N}1$

which ensures that-

- The order in which the elements are being produced along out is identical to the order in which they were consumed.
 - only a limited number of elements can be within the buffer at any time.
- 3) The part of the environment which consumes elements from the buffer can be described as-

$$O_env[X] \equiv [out:seq[X]] \quad DEF(5.4)$$

i.e. an environment which accepts any element from X at any point in time.

- 4) The total system is described through the compound-

$$BB_SYS[X] \equiv [I_env[X]; BB[X]; O_env[X]] \quad DEF(5.5)$$

An *acceptable final state* for the system can be described as-

$$BB_Final[X] \equiv [BB_SYS[X] \mid \#in = Max; in = out] \quad DEF(5.6)$$

i.e. a state where

- all input submitted to the buffer from the environment is consumed by the buffer.
- all input consumed by the buffer has appeared as output from the buffer.

BB_SYS and BB_Final constitute a description of a buffer system where all information except the *interface* between the environment and the buffer is *hidden*. I_env and O_env are components of BB_SYS and document the assumptions made about the environment. BB describes the decisions taken with respect to the component we are about to implement. BB_Final describes the acceptable termination-states of that system.

The *consequences* of the decisions taken so far are:

- 1) The occurrence of a (visible) communication along the in connection can be described as-

$$\begin{aligned} \text{input}[X] = & \\ & (\lambda x : X) \\ & (\lambda BB_SYS)(\mu BB_SYS')(in' = in * \langle x \rangle; out' = out) \end{aligned} \quad DEF(5.7)$$

NB. compare this definition with DEF(4.21).

A communication along the out connection can be described as-

$$\begin{aligned} \text{output}[X] = & \\ & (\lambda x : X) \\ & (\lambda BB_SYS)(\mu BB_SYS')(in' = in; out' = out * \langle x \rangle) \end{aligned} \quad DEF(5.8)$$

The domains of these operations can be derived from the axioms of the BB description.

Example:

The element x can be observed to be consumed by a system b_sys only if

$$\begin{aligned} & b_sys : \{BB_SYS\} \\ \text{and} & \\ & b_sys \in \text{dom}(\text{input}(x)) \end{aligned}$$

We have -

- a) An input can occur whenever we are in a state belonging to

$$\begin{aligned} U\{ \text{dom}(\text{input}(x)) \mid x : X \} = & \\ \{ BB_SYS \mid \#in + 1 \leq Max; (\#in + 1) - \#out \leq B \} = & \\ \{ BB_SYS \mid \#in < Max; \#in - \#out < B \} & \end{aligned}$$

and

- b) An output can occur whenever we are in a state belonging to

$$\begin{aligned} U\{ \text{dom}(\text{output}(x)) \mid x : X \} = & \\ \{ BB_SYS \mid (\exists x : X)(out * \langle x \rangle \subseteq in) \} = & \\ \{ BB_SYS \mid out \neq in \} & \end{aligned}$$

- 2) The system *terminates* because we can present a decreasing variant function

$$\begin{aligned} \text{Var}_{BB} : \{BB_SYS\} \rightarrow \mathbb{N} & \\ \text{Var}_{BB} = (\lambda BB_SYS)(2 * Max - \#in - \#out) & \quad DEF(5.9) \end{aligned}$$

3) The system is *deadlock free*.

The system will terminate in a state which can be described as-

Terminal =

$$\{BB_SYS\} - \bigcup \{ \text{dom}(\text{input}(x)) \mid x:X \}$$

$$- \bigcup \{ \text{dom}(\text{output}(x)) \mid x:X \}$$

=

$$\{BB_SYS\} - \{BB_SYS \mid \text{in} < \text{Max}; \# \text{in} - \# \text{out} < B \}$$

$$- \{BB_SYS \mid \text{out} \neq \text{in} \}$$

and the system is deadlock free if

Terminal \subseteq Final

which is a direct consequence of the following theorem-

$\begin{array}{l} \text{X} \\ \text{THEOREM (5.10)} \\ \text{BB_SYS[X]} \vdash \\ (\text{in} > \text{Max} \wedge \text{out} = \text{in}) \vee (\# \text{in} - \# \text{out} \geq B \wedge \text{out} = \text{in}) \\ \Rightarrow \\ (\# \text{in} = \text{Max} \wedge \text{in} = \text{out}) \end{array}$

(the proof relies on B being *non-zero*)

This completes the *consistency-check* for the decisions taken at the top level of design for a bounded buffer.

NB.

The systems we are concerned with will normally have an additional important property - they are *restartable*. If input is temporarily *blocked* (no requests are submitted to the system for some time) then the system will *deadlock* in a state in which new input must be accepted. Furthermore in such a state all received input must have been output (i.e. all requests have been fulfilled). This property can be formally stated, for BB_SYS, as

$$\{BB_SYS\} - \text{dom}(\bigcup \{ \text{output}(x) \mid x:X \}) \subseteq$$

$$\text{dom}(\bigcup \{ \text{input}(x) \mid x:X \}) \cap \{BB_SYS \mid \text{in} = \text{out}\}$$

For BB_SYS the *restartable* theorem is-

$\begin{array}{l} \text{X} \\ \text{THEOREM (5.11)} \\ \text{Restartable} \\ \text{BB_SYS[X]} \vdash \\ \text{in} = \text{out} \Rightarrow (\# \text{in} - \# \text{out} < B) \wedge (\text{in} = \text{out}) \end{array}$

The approach for *specifying* distributed systems, and for checking whether the decisions taken are *consistent* can be summarised as follows:

- 1) First the designer should decide at which level of detail he wants to present the system. The system is then *partitioned* by the choice of *interfaces* across which we may have *connections*.
- 2) Each part of the system is then described in terms of invariant conditions on the history of the communication along a chosen set of connections.
- 3) The *acceptable final states* for the system are described using predicates on the histories of communications.
- 4) *Termination* is proved by presenting a *variant function*.
- 5) The *consistency-check* is done by either
 - a) proving a deadlock-free theorem, or
 - b) proving a restartable theorem.

5.2. MISCELLANEOUS SYSTEMS.

In order to illustrate the use of the assertional technique for describing behaviour a collection of specifications for some simple, well-known process-components is given in the following sections.

5.2.1. Buffers.

We weaken in turn each of the assertions given for the bounded buffer described in the previous chapter, and get-

- a) An unbounded fifo.

$$\text{FIFO[X]} \equiv [\text{in}, \text{out}; \text{seq[X]} \mid \text{out} \subseteq \text{in}] \quad \text{DEF (5.12)}$$

By placing this component in the environment proposed for the bounded buffer of section 5.1, we get-

$$\text{FIFO_SYS} \equiv [\text{I_env}; \text{FIFO}; \text{O_env}] \quad \text{DEF (5.13)}$$

With the following definition of *acceptable final states*-

$$\text{FIFO_FINAL} \equiv [\text{FIFO_SYS} \mid \# \text{in} = \text{Max}; \text{in} = \text{out}] \quad \text{DEF (5.14)}$$

we can prove that the compound system FIFO_SYS will terminate and is deadlock-free.

In terms of final results the two systems BB_SYS and FIFO_SYS are identical, i.e.

$$\{BB_FINAL\} = \{FIFO_FINAL\}$$

However the two systems are not identical as they have been specified to have different dynamic behaviour. E.g. in a state where

$$\begin{aligned} \#in - \#out &= B \\ (\text{B is the bound for BB_SYS}) \end{aligned}$$

then a BB_SYS will *reject* input until an output has occurred, while the FIFO_SYS can *accept* an input.

b) An Unordered Buffer.

$$\begin{aligned} \text{Unordered_Buffer}[X] &\equiv \\ [in, out: seq[X] \mid (\forall x: X) \text{card}(out \downarrow \{x\}) \leq \text{card}(in \downarrow \{x\})] \\ &\text{DEF (5.15)} \end{aligned}$$

An Unordered_Buffer behaves like an unbounded buffer along the in connection. However while the output produced by the FIFO_SYS at any point in time can be determined to be-

$$in(\#out + 1)$$

the Unordered_Buffer component can produce any element x for which $\text{card}(out \downarrow \{x\}) < \text{card}(in \downarrow \{x\})$

c) A Bounded Unordered Buffer-

$$\begin{aligned} \text{Bounded_Unordered_Buffer}[X] &= \\ [\text{Unordered_Buffer} \mid \#in - \#out \leq B] &\text{DEF (5.16)} \end{aligned}$$

for some $B : \mathbb{N}$

This buffer behaves like a bounded buffer along the in connection, and like an unordered buffer along the out connection.

5.2.2. Queues.

Jobs within an operating system are given unique identification codes in order for the different parts of the system (schedulers for controlling the sharing of resources) to distinguish between those jobs. The histories of the *arrival* and the *departure* to and from the queues of jobs (wait-queues, ready-queues, spooling-queues) within an operating system can therefore be described as-

$$\text{id_seq}[X] = \{ s: seq[X] \mid s^{-1} \in X \rightarrow \mathbb{N} \}$$

Hence for any connection or channel in such a system,

$$\begin{aligned} \text{chan} &: \text{id_seq}[X] \\ \text{the function defined as} \\ \text{rel_time} &: X \rightarrow \mathbb{N} \\ \text{rel_time} &= \text{chan}^{-1} \end{aligned}$$

gives the relative time of the occurrence of communications along chan.

The id_seq will in the following be used to describe the behaviour of various queues. Note that the following descriptions do not include a description of the environment, thus a complete discussion of the behaviour of the component presented cannot be given.

a) Unordered Queue.

$$\begin{aligned} \text{QUEUE}[J] &\equiv \\ [\text{arrive, depart: id_seq}[J] \mid \text{ran}(\text{depart}) \subseteq \text{ran}(\text{arrive})] \\ &\text{DEF (5.17)} \end{aligned}$$

which states the obvious property that only jobs which have arrived in a queue can depart from it. Note that the time of departure is not related to the time of arrival. A job within a QUEUE may have its departure delayed unreasonably often, because as long as more than one job is within the queue it is non-deterministic which job actually leaves the queue, hence if the rate of input is high a job can be passed over time after time. By adding a description of priority mechanisms to the QUEUE specification we can impose a fairer scheduling for departures.

b) Ordered Queues.

$$\begin{aligned} \text{FIFO_QUEUE}[J] &\equiv \\ [\text{QUEUE}[J] \mid \text{depart} \subseteq \text{arrive}] &\text{DEF (5.18)} \end{aligned}$$

In a correct implementation of a FIFO_QUEUE the jobs will depart in the order of arrival. This scheduling mechanism may however have undesirable run-time characteristics. Unnecessary delays may occur if the process component which receives jobs along the departure connection is temporarily unwilling to accept the next job, but could have serviced another queued job.

An alternative design would be-

c) Limited-overtaking (or 'fair') queue

$$\begin{aligned} \text{FAIR_QUEUE} &\text{---} J \text{---} \text{DEF (5.19)} \\ \text{QUEUE}[J] & \\ \text{---} & \\ l + \# \text{depart} - \min(\text{arrival_times}) &\leq L \\ \text{where} & \\ \text{arrival_times} &= \text{dom}(\text{arrive} \downarrow \text{queued}); \\ \text{queued} &= \text{ran}(\text{arrive}) - \text{ran}(\text{depart}) \end{aligned}$$

(for some limit L which is non-zero)

The FAIR_QUEUE describes a queue which ensures that no jobs can be overtaken by more than L other jobs.

Example:

let

```
L = 3;
arrive = <j1 j2 j3 j4 j5 j6>
depart = <j2 j4 j3>
```

then

$$l+\#depart - \min(\text{arrival_times}) = 1+3 - 1 = 3$$

and the departure of a job other than the job j1 would invalidate the axioms given. If j1 departs we get-

```
depart = <j2 j4 j3 j1>
```

and

$$l+\#depart - \min(\text{arrivaltimes}) = 1+4 - 5 = 0$$

5.3. CRITICAL SECTIONS AND MUTUAL EXCLUSION.

Consider the mutual exclusion problem. This is a problem of devising a scheduling mechanism for concurrently operating jobs or processes which ensures that a job does not *enter* a *critical section* where it makes use of a shared resource if other jobs, currently within their critical section, have used up the resource.

Let us first formalise *Mutual Exclusion*.

Let P be a set of processes. The mutual exclusion can be described in terms of a relation over P , i.e.

Mutex $\frac{\quad}{P}$ DEF (5.20)

exclude : $P \leftrightarrow P$

$\text{id}(P) \cap \text{exclude} = \{\}$;
 $\text{exclude}^{-1} = \text{exclude}$

which states that-

- nobody can exclude themselves (the Mutex relation is *irreflexive*).
- the exclusion is mutual (the Mutex relation is *symmetric*).

Examples:

- Let us describe a small readers/writers system with three readers and one writer.

$$RW_system = \{ r_1, r_2, r_3, w_1 \}$$

The readers exclude the writer and the writer excludes all readers, i.e.

RW : RW_system \leftrightarrow RW_system

$$RW = \{ r_1 \leftrightarrow w, r_2 \leftrightarrow w, r_3 \leftrightarrow w, \\ w \leftrightarrow r_1, w \leftrightarrow r_2, w \leftrightarrow r_3 \}$$

DEF (5.21)

and

RW \in {Mutex[RW-System]}

- Let

$$Phil = \{ p_1, p_2, p_3, p_4, p_5 \}$$

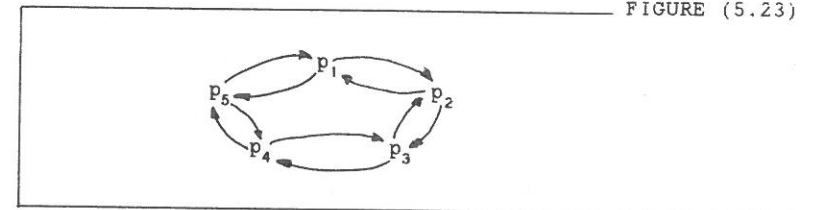
and

PH : Phil \leftrightarrow Phil

$$PH = \{ p_1 \leftrightarrow p_2, p_2 \leftrightarrow p_3, p_3 \leftrightarrow p_4, p_4 \leftrightarrow p_5, \\ p_5 \leftrightarrow p_4, p_4 \leftrightarrow p_3, p_3 \leftrightarrow p_2, p_2 \leftrightarrow p_1 \}$$

DEF (5.22)

which can be illustrated as-



then the relation PH is a Mutex relation, where each individual excludes *neighbours*.

A simple system which monitors a set of jobs which alternately execute outside and inside a critical section, can be characterised in terms of two components:

- The non-critical section can be described as-

ENV $\frac{\quad}{J}$ DEF (5.24)

enter, leave : seq[J]

$$(\forall j:J) (\text{card}(\text{enter} \downarrow \{j\}) - \text{card}(\text{leave} \downarrow \{j\}) \leq 1 \wedge \\ \text{card}(\text{enter} \downarrow \{j\}) \leq \text{MAX}(j))$$

for some $\text{MAX} \in J \rightarrow \mathbb{N}$, which ensures that all jobs only *enter* the critical section a finite number of times.

2) Given a Mutex relation M , then a critical section can be described as-

CS	J	DEF (5.25)
$\text{enter, leave : seq}[J]$ $M(\text{inside}) \cap \text{inside} = \{\};$ $(\forall j:J) (\text{card}(\text{enter}\downarrow\{j\}) > \text{card}(\text{leave}\downarrow\{j\}))$ <p>where</p> $\text{inside} =$ $\{j:J \mid \text{card}(\text{enter}\downarrow\{j\}) > \text{card}(\text{leave}\downarrow\{j\})\}$		

which constrains the behaviour along the *enter* connection so that no job can get inside if another job to which it is related through the mutual exclusion relation, is currently inside. The obvious property that no jobs can leave before they enter is also included in the description.

The final states of the system-

$$\text{SYS} \approx [\text{ENV}; \text{CS}] \quad \text{DEF (5.26)}$$

are the states where all jobs have terminated and are outside the critical section, i.e.

FINAL	J	DEF (5.27)
SYS $(\forall j:J) (\text{card}(\text{enter}\downarrow\{j\}) = \text{MAX}(j) \wedge$ $\text{card}(\text{leave}\downarrow\{j\}) = \text{MAX}(j))$		

The system SYS terminates and is deadlock free.

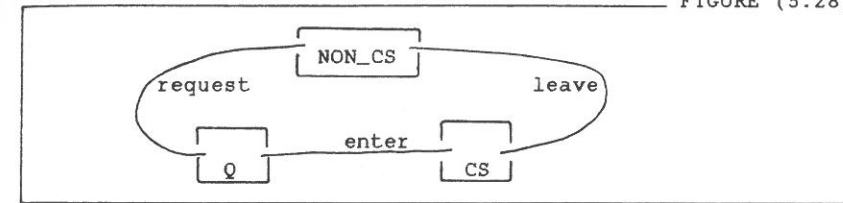
NB, taking the relation PH in DEF(5.22) as the mutex relation we have a top level specification of the *dining philosophers* problem, and taking RW in DEF(5.21) we have a description of a small readers/writers system.

A property which is normally discussed in connection with mutual exclusion (in addition to deadlock properties) is the absence of *starvation*. Starvation may occur when a process (or job) can be excluded from entering the critical section by two or more other independent processes. The two processes can either intentionally co-operate in excluding the process or accidentally be scheduled in such a way that one of the two processes will always be within the critical section.

How can we change our specification in such a way that we impose a scheduling strategy that prevents starvation?

Consider the following system-

FIGURE (5.28)



where NON_CS makes up the environment for the system being designed and is defined as-

$$\text{NON_CS}[P] \approx \text{ENV}[P][\text{request}/\text{enter}] \quad \text{DEF (5.29)}$$

and represents a collection of independent processes (in their non-critical section) which alternately *request-to-enter* and *leave* the critical section.

Q represents a queue of processes which have requested permission to enter the critical section but are temporarily delayed in doing so. The Q can be described as-

- a) An unordered queue (see DEF(5.17))

$$Q[P] \approx \text{QUEUE}[P][\text{request}/\text{arrive}, \text{enter}/\text{depart}] \quad \text{DEF (5.30)}$$

which does not prevent starvation.

- b) A fifo queue (see DEF(5.18))

$$Q[P] \approx \text{FIFO_QUEUE}[P][\text{request}/\text{arrive}, \text{enter}/\text{depart}] \quad \text{DEF (5.31)}$$

which prevents starvation but may have undesirable run-time characteristics.

- c) A limited overtaking queue (see DEF(5.19))

$$Q[P] \approx \text{FAIR_QUEUE}[P][\text{request}/\text{arrive}, \text{enter}/\text{depart}] \quad \text{DEF (5.32)}$$

which prevents starvation and only causes unnecessary delays in situations where a process has been overtaken too often.

The system which is illustrated in FIGURE(5.28) (with the Q being defined as a FAIR_QUEUE) constitutes a top level specification of a system which schedules several independently executing processes in such a way that -

- 1) no two processes which are related through a given Mutex relation are within the critical section at the same time.
- 2) The system is deadlock free, with respect to the following definition of the final state-

FINAL_CS	P	DEF (5.33)
NON_CS[P]; Q[P]; CS[P]		
$(\forall p:P)$ $(\text{card}(\text{request}\downarrow\{p\}) = \text{MAX}(p) \quad \wedge$ $\text{card}(\text{request}\downarrow\{p\}) = \text{card}(\text{enter}\downarrow\{p\}) \quad \wedge$ $\text{card}(\text{enter}\downarrow\{p\}) = \text{card}(\text{leave}\downarrow\{p\}))$		

- 3) No processes are being starved. By introducing the queue into the specification we introduced new operations which enabled us 1) to make a distinction between a request to enter and the actual entering of the critical section, 2) to define starvation and 3) to suggest a 'dynamics' which prevents it.

5.4. PREVENTION OF DEADLOCK.

An inconsistency in the decisions made with respect to - 1) the invariants of the system and 2) the final state - will result in *Deadlock*, i.e. if a system is realised following the rules given for the *behaviour* of the system (the invariant conditions) then this system cannot be guaranteed to reach a *final state*. The importance of the relationship between invariant conditions and the definition of the final state has been discussed in section 3.3.1.

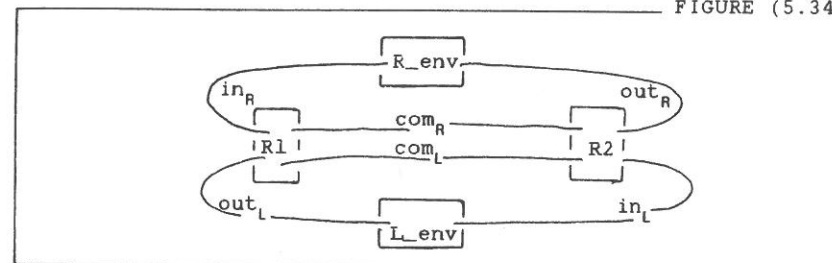
The approach for specifying distributed systems introduced in the preceding chapters enables the designer to investigate deadlock properties solely from the information provided in the specification. If a specification cannot be proved to be deadlock free the designer should either change the *behaviour* of the system by modifying the invariant of the system or change the definition of the *final state*.

We will in this section give an example of a small system which cannot be proved to be deadlock free from its first specification. A modification to the specification which *strengthens* the invariant conditions will ensure that behaviour which may lead to deadlock is avoided.

Consider a distributed system consisting of four independent process components. Two of the components are non-sharable resources (R1, R2), and two are components in the environment. The first environment component (R_env) is in control of starting a sequence of jobs (R_jobs). Each R_job will, when initiated, first occupy resource R1. Next it acquires resource R2. On completion it signals the R_env. The other environment (L_env) controls a sequence of L_jobs. The L_jobs behave like R_jobs except for the order in which they acquire the resources R1 and R2. An L_job requires access to R2 before it requires access to R1.

The system may be illustrated as follows:

FIGURE (5.34)



Given two initial sequences:

```
init_R : seq[R];
init_L : seq[L]
```

we define

$R_env[R] \equiv [in_R, out_R : seq[R] \mid in_R \subseteq init_R]$ DEF (5.35)

$L_env[L] \equiv [in_L, out_L : seq[L] \mid in_L \subseteq init_L]$ DEF (5.36)

i.e. the R_env will transmit all elements from $init_R$ along in_R in the order they appear in $init_R$ and L_env will transmit elements from $init_L$ along in_L .

R1 is described as-

R1 R L DEF (5.37)

```
in_R, com_R : seq[R];
com_L, out_L : seq[L]
```

$com_R \subseteq in_R; out_L \subseteq com_L;$

$\delta l_R + \delta l_L \leq 1$

where

$\delta l_R = \#in_R - \#com_R$

$\delta l_L = \#com_L - \#out_L$

The first two axioms state that the resource-controller R1 will transmit the jobs received along in_R to the channel com_R and the jobs received along com_L to the out_L channel. The third axiom ensures that only a single job can use the resource at any point in time.

Similarly for R2 we have-

R2	R L	DEF (5.38)
$\begin{aligned} in_L, com_L &: seq[L]; \\ com_R, out_R &: seq[R] \end{aligned}$		
$\begin{aligned} com_L \subseteq in_L; out_R \subseteq com_R; \\ \delta 2_R + \delta 2_L \leq 1 \end{aligned}$		
<p>where</p> $\begin{aligned} \delta 2_L &= \#in_L - \#com_L \\ \delta 2_R &= \#com_R - \#out_R \end{aligned}$		

In a final state for the system-

RL_SYS	R L	DEF (5.39)
$\begin{aligned} R_env[R,L]; \\ L_env[R,L]; \\ R1[R,L]; \\ R2[R,L] \end{aligned}$		

we expect all jobs originally in the input queues ($init_L$ and $init_R$) to have completed their execution, i.e.

RL_FINAL	R L	DEF (5.40)
$RL_SYS[R,L]$		
$\begin{aligned} in_R = init_R; in_L = init_L; \\ in_R = out_R; in_L = out_L \end{aligned}$		

NB. if we are in a final state then all resources are free, i.e.

$$RL_FINAL \vdash \delta 1_R = \delta 1_L = \delta 2_R = \delta 2_L = 0$$

We now observe-

- 1) The specified system will terminate. A variant function can be defined as -

$$\begin{aligned} Var = (\lambda RL_SYS) \\ (\#init_R * 3 - \#in_R - \#com_R - \#out_R + \\ \#init_L * 3 - \#in_L - \#com_L - \#out_L) \end{aligned} \quad \text{DEF (5.41)}$$

- 2) The terminal states for the system can be described as -

RL_TERM	R L	DEF (5.42)
$RL_SYS[R,L]$		
$(in_R = init_R \vee \delta 1_R = 1 \vee \delta 1_L = 1); \quad (1)$		
$(\delta 1_R = 0 \vee \delta 2_R = 1 \vee \delta 2_L = 1); \quad (2)$		
$(\delta 1_R = 0); \quad (3)$		
$(in_L = init_L \vee \delta 2_R = 1 \vee \delta 2_L = 1); \quad (4)$		
$(\delta 2_L = 0 \vee \delta 1_R = 1 \vee \delta 1_L = 1); \quad (5)$		
$(\delta 2_L = 0) \quad (6)$		

NB.

- (1) denotes the states where activity along in_R are blocked.
 (2) denotes the states where activity along com_R are blocked.
 (3) denotes the states where activity along out_R are blocked.
 (4,5,6) denote similar conditions for L -jobs.

- 3) The system cannot be guaranteed to reach a final state as we cannot prove-
 $RL_TERM[R,L] \vdash RL_FINAL[R,L]$

However we can prove-

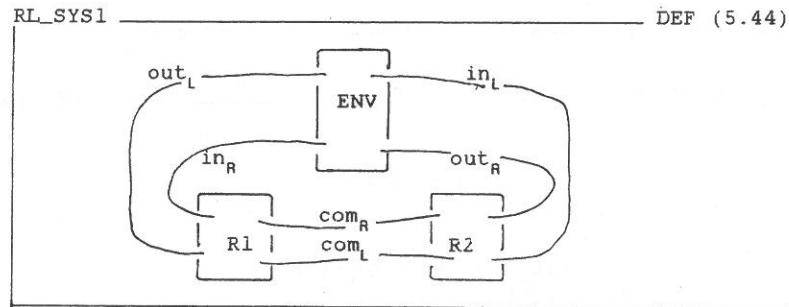
RL_TERM	R L	THEOREM (5.43)
$RL_TERM[R,L] \vdash$		
$(in_R = init_R \wedge in_L = init_L \wedge in_R = out_R \wedge in_L = out_L)$		
\vee		
$(\delta 1_R = 1 \wedge \delta 2_L = 1)$		

which states that a RL_SYS will terminate in either an acceptable final state or in a state where an R -job occupies $R1$ (waiting for $R2$) and an L -job occupies $R2$ (waiting for $R1$).

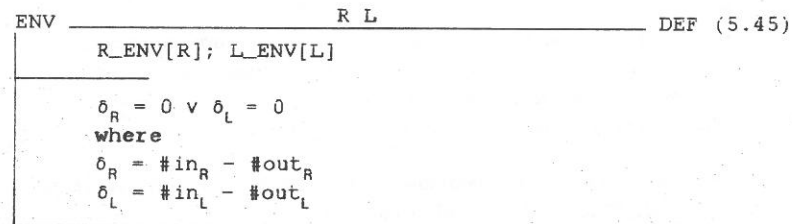
There are several ways of modifying the specification in order to eliminate undesirable deadlocks (see section 3.3.1.). For this particular system it seems reasonable to maintain the given specification of the final state and solve the 'deadlock' by strengthening the invariant conditions thereby modifying the dynamics of the system.

Two solutions will be given in the following:

- 1) The deadlock can be avoided by simply making further assumptions about the environment. If we can rely on the environment never to start an R -job while an L -job is in progress and vice versa, then the system will be deadlock-free. Such a system may be illustrated as-



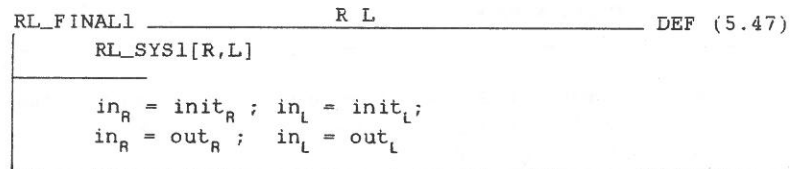
where R1 and R2 are unchanged and-



The new system-

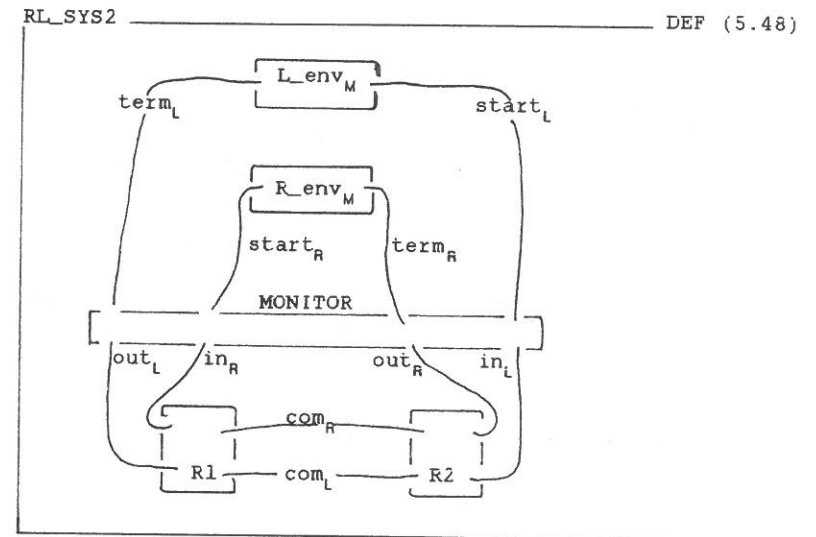


with the final state-

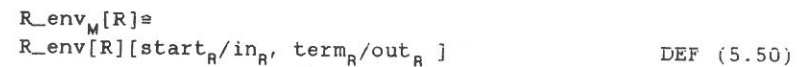
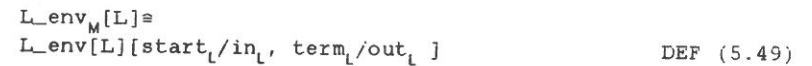


terminates and can be proved to be deadlock free.

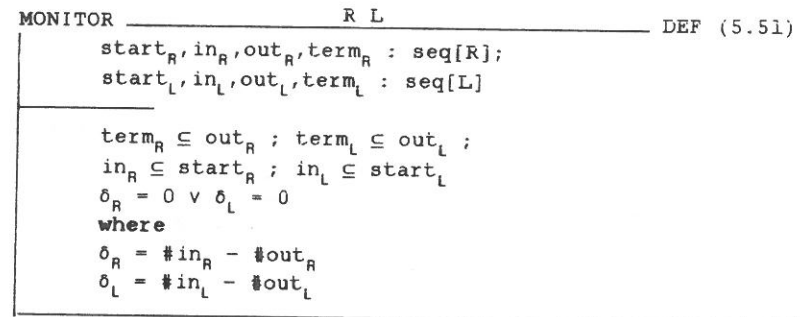
- 2) The behaviour of the system can also be changed by adding a *monitor* which coordinates the L-traffic and the R-traffic. Such a system can be illustrated as-



where R1 and R2 are unchanged, and the environments are changed to communicate with the introduced monitor, i.e.



The MONITOR is a buffer between the resource controllers and the environment and delays R_jobs if an L_job is in progress and vice versa.



RL_SYS2 can with an obvious definition for the final state be proved to be deadlock-free.

Note that we used the result from the analysis of the first system (RL_SYS) to guide design modifications. However the process of removing deadlock from a design cannot in general be automated. The presence of deadlock reveals that the designer has not fully understood the consequences of his design. His design is *inconsistent* and must be modified, but exactly what must be changed is solely his decision.

5.5. A COMMUNICATION NETWORK.

In this section a specification of a communication network will be developed. The design decisions taken will be formally stated and motivated. A formal description of the behaviour of a single station (node) will be given, as well as a description of the static properties of the network as a whole. These descriptions will be used to verify that the *global behaviour* of a network, which consists of stations whose *local behaviour* complies with the given specification, agrees with some independently stated *rules* for networks.

5.5.1 A Discussion of a Communication Network.

A NETWORK consists of a collection of STATIONS. In this formal description the set of stations will be denoted by-

$$\text{St}$$

The value which is communicated over the NETWORK is called a PACKET. Packets are denoted by-

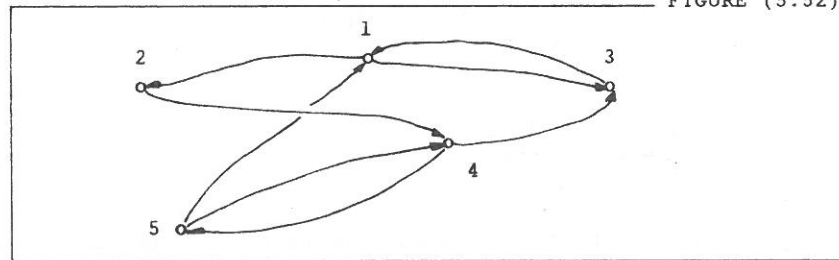
$$\text{Pk}$$

All packets have a SOURCE station and a DESTINATION station-

$$\text{source, dest} : \text{Pk} \rightarrow \text{St}$$

Example:

Consider the network illustrated in Figure (5.52) consisting of five stations.

$$\{ 1, 2, 3, 4, 5 \}$$


Given a set of packets.

$$\{ a, b, c \}$$

We might have.

$$\text{source} = \{ (a,1), (b,5), (c,2) \}$$

$$\text{dest} = \{ (a,5), (b,3), (c,5) \}$$

The set of packets which originates in station i is

$$\text{source}^{-1}(\{i\})$$

The set of packets which can be consumed by i is

$$\text{dest}^{-1}(\{i\})$$

For the example above we get.

$$\text{source}^{-1}(\{1\}) = \{a\}$$

$$\text{dest}^{-1}(\{5\}) = \{a, c\}$$

Let us first formalise some obvious properties for stations in a network. Stations either *consume* or *transmit* packets. In order to formalise this statement we introduce the *history* of the input performed by each station-

$$\text{In} : \text{St} \rightarrow \text{seq}[\text{Pk}]$$

and a *history* of the output performed by each station-

$$\text{Out} : \text{St} \rightarrow \text{seq}[\text{Pk}]$$

The output from station i , $\text{Out}(i)$, will be written Out_i ; similarly for In .

We can now state some requirements for the behaviour of a station.

a) Only packets originating in (or previously received by) a station can be transmitted along the outgoing connections for that station. We have-

$$(\forall i:\text{St})(\text{ran}(\text{Out}_i) \subseteq \text{source}^{-1}(\{i\}) \cup \text{ran}(\text{In}_i))$$

b) A packet arriving at its destination must be consumed- i.e. it must not be retransmitted, therefore we strengthen a)

$$(\forall i:\text{St})$$

$$(\text{ran}(\text{Out}_i) \subseteq (\text{source}^{-1}(\{i\}) \cup \text{ran}(\text{In}_i)) - \text{dest}^{-1}(\{i\}))$$

DEF (5.53)

c) If we assume that packets are distinguishable- the fact that no packet can be transmitted from a particular station more than once can be formalised-

$$(\forall i:\text{St})(\text{Out}_i^{-1} \in \text{Pk} \rightarrow \mathbf{N})$$

DEF (5.54)

b) and c) record the decisions taken with respect to the behaviour of a station, ignoring the fact it may have several *input channels* and several *output channels*. We may say that a station's behaviour is like the behaviour of an *unlimited unordered buffer*, i.e. a station will always accept incoming packets, but these packets may be stored for later transmission. Such a system is called a *store and forward system*.

Alternative designs will be discussed below:

1) Limited buffering capacity.

In order to describe stations with limited buffering capacity we might have added the in variant, (for some limit L).

$$(\forall i:St) (\text{card}(\text{ran}(In_i) - \text{dest}^{-1}(\{i\})) - \text{card}(Out_i) \leq L)$$

However, a limit on the buffer-store for a station in a network may cause the network to deadlock.

2) Scheduling.

No policy for scheduling the outgoing packets has been suggested. Hence we allow for the possibility of *unfair scheduling* where a packet within a station may be *starved* (overtaken by other packets an unreasonable number of times).

A *FIFO* scheduling might be described as follows:

For a *relay-station* i where

$$\text{dest}^{-1}(\{i\}) = \text{source}^{-1}(\{i\}) = \{\}$$

a *FIFO* ordering could be imposed as usual by requiring that the history sequence describing the output is a *prefix* of the sequence describing the input, i.e.

$$Out_i \subseteq In_i$$

For a *terminal-station* i , where

$$\text{dest}^{-1}(\{i\}) \neq \{\} \vee \text{source}^{-1}(\{i\}) \neq \{\}$$

the *ordering* axiom could be expressed as follows:

$$\text{ignore}(\text{source}^{-1}(\{i\}))(Out_i) \subseteq \text{ignore}(\text{dest}^{-1}(\{i\}))(In_i)$$

where

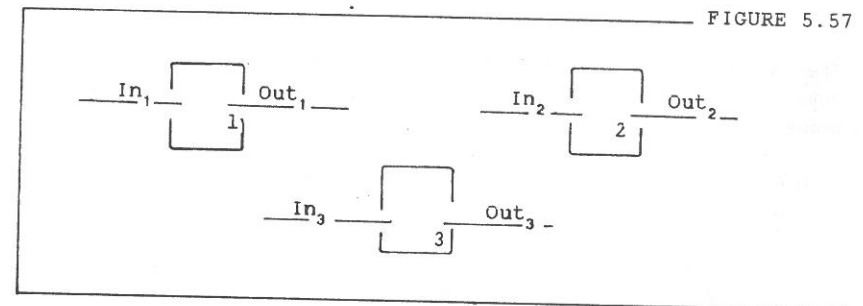
$$\begin{array}{l} \text{X} \\ \text{DEF (5.56)} \\ \text{ignore}[X] : \mathcal{P}(X) \rightarrow \text{seq}[X] \rightarrow \text{seq}[X] \\ (\forall S:\mathcal{P}(X); s:\text{seq}[X]) \\ \text{ran}(\text{ignore}(S)(s)) = \text{ran}(s) - S \wedge \\ (\exists f:\text{Monotone}) (\text{ignore}(S)(s) = s \circ f) \\ \text{where} \\ \text{Monotone} = \{f:\mathbb{N} \rightarrow \mathbb{N} \mid (\forall x,y:\mathbb{N}) x \leq y \Rightarrow f(x) \leq f(y)\} \end{array}$$

NB *ignore* removes elements from a sequence. The result is a new sequence where the order of the remaining elements is the same as the ordering between these elements in the original sequence, e.g.

$$\text{ignore}(\{a,c\})(\langle a b c d e \rangle) = \langle b d e \rangle$$

A *FIFO* scheduling policy could cause unnecessary runtime delays, because the station to which a scheduled packet is to be transmitted could be temporarily engaged in other communications. For these reasons we will not impose a *FIFO* scheduling policy.

DEF(5.53) and DEF(5.54) gave some dynamic properties for a station, without any reference to the topology of the network, i.e. so far we have only described a collection of dis-connected stations-



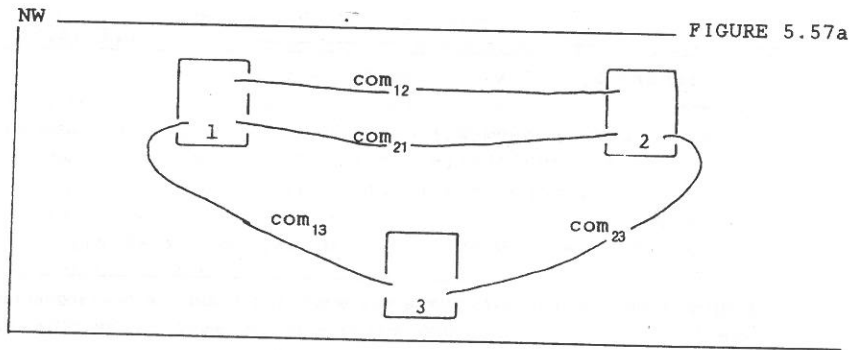
The connections between the individual stations are described by:

$$\text{Network} : St \leftrightarrow St$$

where we assume that the number of connections is finite, i.e.

$$\text{Network} \in \mathcal{F}(St \times St)$$

hence we might get:



When packets are sent along the connections of the network, the information which changes is:

Communication_between : $St \times St \rightarrow seq[Pk]$ where
 $dom(Communication_between) = Network$

that is, when a packet is sent from station i to station j , the *history sequence* $Communication_between(i, j)$ will be extended.

Let com_{ij} denote (as in FIGURE(5.57)) the

Communication_between(i, j)

The correct behaviour of a station depends on its capability to transmit each packet along a proper connection. We will in what follows discuss two possibilities for proper routing of packets.

- a) A packet should not be sent to a station from which the packet's destination cannot be reached, i.e.

$(\forall i, j: St \mid (i, j) \in Network)$
 $(\forall p: Pk \mid p \in ran(com_{ij})) (j, dest(p)) \in Network^*$

NB $Network^*$ is a description of the stations which are connected through the network.

The axiom above does not prevent a packet from 'entering a loop' or 'traveling' erratically round the network without ever reaching its destination. We therefore strengthen the axiom:

$(\forall i, j: St \mid (i, j) \in Network)$
 $(\forall p: Pk \mid p \in ran(com_{ij}))$
 $((j, dest(p)) \in Network^* \wedge$
 $distance(Network)(j, dest(p)) <$
 $distance(Network)(i, dest(p)))$
 where

X

DEF (5.58)

distance : $(X \leftrightarrow X) \rightarrow (X \times X \rightarrow N)$

$(\forall r : X \leftrightarrow X)$
 $(\forall i, j: X \mid (i, j) \in r^*)$
 $(distance(r)(i, j) = \min\{n: N \mid (i, j) \in r^n\})$

This condition insures that all packets visit a finite number of stations, and travel along a route of minimal length.

- b) Condition a) does not allow a station to send a packet along a connection which is not on the shortest route for that packet. The following specification allows a station to choose a longer but possibly faster route. It is still guaranteed that all packets will reach their destinations in a finite number of steps.

We first formalise the concept of a proper route for a packet, which is a route within the network where no station is visited more than once, hence

X

DEF (5.59)

proper_routes : $(X \leftrightarrow X) \rightarrow \mathcal{P}(seq[X])$

$(\forall r : X \leftrightarrow X)$
 $proper_routes(r) =$
 $\{rt: seq[X] \mid next(rt) \subseteq r; rt^{-1} \in X \rightarrow N\}$

NB next is defined in section 2.3.2.

for the network illustrated in FIGURE(5.57a) we have:

$\langle 1\ 2\ 3 \rangle \in proper_routes(NW)$
 $\langle 1\ 3 \rangle \in proper_routes(NW)$
 $\langle 3\ 1 \rangle \notin proper_routes(NW)$
 $\langle 1\ 2\ 1\ 3 \rangle \notin proper_routes(NW)$

The requirements for a packet router for the network.

Packet_router : $Pk \rightarrow seq[St]$

can now be formalised:

All routes are proper routes:

$ran(Packet_router) \subseteq proper_routes[St](Network)$

Exactly the packets which can reach their destination from their source can be given a route by the router:

$$\text{dom}(\text{Packet_router}) = \{p:\text{Pk} \mid (\text{source}(p), \text{dest}(p)) \in \text{Network}^*\}$$

A route for a packet 'starts' at its source and 'ends' at its destination:

$$\begin{aligned} \text{dest} &= \text{last} \circ \text{Packet_router} \\ \text{source} &= \text{first} \circ \text{Packet_router} \end{aligned}$$

NB. separate packets sent between the same source and destination may take (or be given) different routes.

The use of output channels is determined by the packet router in the following way:

$$\begin{aligned} (\forall i : \text{St}) \\ (\forall j : \text{St} \mid (i, j) \in \text{Network}) \\ (\forall p : \text{Pk} \mid p \in \text{ran}(\text{com}_{ij})) \\ (i, j) \in \text{next}(\text{Packet_router}(p)) \end{aligned}$$

5.5.2. A Specification of a Network.

The decisions taken in the previous section will be summarised in this section.

The *static* properties of the network presented are described in DEF(5.60). A description of source and destination (introduced above) is no longer necessary as the source and destination of packets are the *first* and the *last* stations along their route. Packets are all the packets ever processed by the Network, hence we must require that at least these packets have a route within the network. *Start_in* gives for each station the set of packets originating in that station. Previously we used the *inverse* of the source function to describe these packets.

St Pk	DEF (5.60)
Network	: $\mathcal{F}(\text{St} \times \text{St})$;
Packet_router	: $\text{Pk} \rightarrow \text{seq}[\text{St}]$;
Packets	: $\mathcal{P}(\text{Pk})$;
Start_in	: $\text{St} \rightarrow \mathcal{P}(\text{Pk})$
$\begin{aligned} \text{ran}(\text{Packet_router}) &\subseteq \text{proper_routes}(\text{Network}); \\ \text{Packets} &\subseteq \text{dom}(\text{Packet_router}); \\ \text{Start_in} &= (\text{first} \circ \text{Packet_router})^{-1} \downarrow \text{Packets} \end{aligned}$	

The *dynamic* properties are described in DEF(5.61). The properties are described in terms of *invariant conditions* on the histories of the communications between

the stations. Axioms (2) and (3) are substituted for the axioms DEF(5.53) and DEF(5.54) respectively; they have been simplified as the Packet_router, axiom (4), insures that-

- a) a packet cannot 'leave' its destination,
- b) a packet can 'leave' a station along only one outgoing connection.

NETWORK	St Pk	DEF (5.61)
$\text{com} : (\text{St} \times \text{St}) \rightarrow \text{seq}[\text{Pk}]$		
$\text{dom}(\text{com}) = \text{Network};$		
$(\forall i : \text{St})$		
$(\forall j : \text{St} \mid (i, j) \in \text{Network})$		
$(\text{ran}(\text{com}(i, j)) \subseteq \text{Start_in}(i) \cup \text{IN}(i) \wedge$		
$\text{com}(i, j)^{-1} \in \text{Pk} \rightarrow \text{N} \wedge$		
$((\forall p : \text{Pk} \mid p \in \text{ran}(\text{com}(i, j)))$		
$(i, j) \in \text{next}(\text{Packet_router}(p)))$		
$\text{where } \text{IN}(i) = \cup\{\text{ran}(\text{com}(k, i) \mid k : \text{Network}^{-1}(\{i\})\}$		

The Initial state of a network system is a system where no communication has taken place-

INITIAL_NW	St Pk	DEF (5.62)
$\text{NETWORK}[\text{St}, \text{Pk}]$		
$(\forall (i, j) : \text{St} \times \text{St} \mid (i, j) \in \text{Network}) (\text{com}(i, j) = \langle \rangle)$		

The final state of a network system is a system where all packets have reached their destination-

FINAL_NW	St Pk	DEF (5.63)
$\text{NETWORK}[\text{St}, \text{Pk}]$		
$(\forall i : \text{St} \mid i \in \text{dom}(\text{Network}))$		
$(\forall p : \text{Pk} \mid p \in (\text{Start_in}(i) \cup \text{IN}(i)) - \text{OUT}(i))$		
$(i = \text{last}(\text{Packet_router}(p)))$		
$\text{where } \text{IN}(i) = \cup\{\text{ran}(\text{com}(k, i) \mid k : \text{Network}^{-1}(\{i\})\}$		
$\text{and } \text{OUT}(i) = \cup\{\text{ran}(\text{com}(i, j) \mid j : \text{Network}(\{i\})\}$		

5.5.3. An Analysis of the Network.

The given network will be analysed with respect to *termination* and *deadlock*. Under the assumption that a finite number of packets are being 'submitted' to the network system, i.e.

$$\text{Packets} \in \mathcal{F}(\text{Pk}),$$

we will verify that

- 1) The operations of the system cease after the occurrence of a finite number of communications (*termination*).
- 2) The network system only terminates in an acceptable final state - as described in FINAL-NW (*deadlock-free*).

Termination.

Termination is guaranteed if we can present a *decreasing variant function*. A communication between station i and station j can be described as-

Observations $\frac{\text{Pk}}{\text{DEF (5.64)}}$

$\text{OBS}_{ij} : \text{Pk} \rightarrow \{\text{NETWORK}\} \rightarrow \{\text{NETWORK}\}$
$\text{OBS}_{ij} = (\lambda p:\text{Pk}) (\lambda \text{NETWORK}) (\mu \text{NETWORK}') \\ (\text{com}' = \text{com} \oplus \{(i,j) \rightarrow \text{com}(i,j) * \langle p \rangle\})$

NB. (\oplus) is the *function overriding operator*.

According to the definition above a communication will *increase* the length of the *history* of the communications along a single connection. Therefore the following function must be an *increasing variant function* -

$$\begin{aligned} \text{V1} &: \{\text{NETWORK}\} \rightarrow \mathbb{N} \\ \text{V1} &= (\lambda \text{NETWORK}) (\text{Sigma}(\text{com}, \text{card} \circ \text{ran})) \end{aligned}$$

where Sigma is defined as

$\frac{\text{X}}{\text{DEF (5.65)}}$

$\text{Sigma} : (\mathcal{F}(X) \times (X \rightarrow \mathbb{N})) \rightarrow \mathbb{N}$
$(\forall s:\mathcal{F}(X); f:X \rightarrow \mathbb{N} \mid s \subseteq \text{dom}(f); s \neq \{\}) \\ \text{Sigma}(s,f) = f(\tau(s)) + \text{Sigma}(s - \{\tau(s)\}, f) \\ \text{Sigma}(\{\}, f) = 0$

According to axiom (3) of NETWORK a packet can appear only once along any connection, hence

$$(\forall \text{nw}:\text{NETWORK}) \text{V1}(\text{nw}) \leq \text{card}(\text{Packets}) * \text{card}(\text{Network})$$

Therefore the following function is a *decreasing variant function*-

$$\begin{aligned} \text{V} &: \{\text{NETWORK}\} \rightarrow \mathbb{N} \\ \text{V} &= (\lambda \text{nw}:\text{NETWORK}) \\ &\quad ((\text{card}(\text{Packets}) * \text{card}(\text{Network})) - \text{V1}(\text{nw})) \end{aligned}$$

NB, in a termination state 's' we do not guarantee that $\text{V}(s) = 0$.

Deadlock.

We can describe the activity along channel (i,j) as follows:

$\frac{\text{Pk}}{\text{DEF (5.66)}}$

$\text{activity_along}_{ij} : \{\text{NETWORK}\} \leftrightarrow \{\text{NETWORK}\}$
$\text{activity_along}_{ij} = \cup \{ \text{OBS}_{ij}(p) \mid p : \text{Pk} \}$

A termination state of a network is a state where no activity can take place, and can be described as:

$\frac{\text{DEF (5.67)}}$

$\text{Terminal} : \mathcal{P}(\{\text{NETWORK}\})$
$\text{Terminal} = \\ \{\text{NETWORK}\} - \cup \{ \text{dom}(\text{activity_along}_{ij}) \mid (i,j) : \text{Network} \}$

A system is *deadlock-free* if it only terminates in acceptable final states.

For this network system we have:

$\frac{\text{THEOREM (5.68)}}$

$\vdash \text{Terminal} \subseteq \{\text{FINAL_NW}\}$

The proof of THEOREM(5.68) is directly derivable from LEMMA(5.70) and LEMMA(5.71)

We first describe the states of the network in which a communication is expected to take place.

$\frac{\text{Pk}}{\text{LEMMA (5.69)}}$

$\vdash \text{dom}(\text{activity_along}_{ij}) = \\ \{\text{NETWORK} \mid (\exists p:\text{Pk}) \\ (p \in (\text{Start_in}(i) \cup \text{IN}(i)) - \text{OUT}(i)) \wedge \\ (i,j) \in \text{next}(\text{Packet_router}(p)) \}$
$\vdash \cup \{ \text{dom}(\text{activity_along}_{ij}) \mid (i,j) : \text{Network} \} = \\ \{\text{NETWORK} \mid (\exists p:\text{Pk}) (\exists (i,j):\text{Network}) \\ (p \in (\text{Start_in}(i) \cup \text{IN}(i)) - \text{OUT}(i)) \wedge \\ (i,j) \in \text{next}(\text{Packet_router}(p)) \}$

NB the functions IN and OUT are defined as in FINAL_NW DEF(5.63).

LEMMA(5.70) describes the termination states:

Pk	LEMMA (5.70)
$\vdash \text{Terminal} =$ $\{ \text{NETWORK} \mid (\forall (i,j):\text{Network}) (\forall p:\text{Pk})$ $(p \notin (\text{Start_in}(i) \cup \text{IN}(i)) - \text{OUT}(i)) \vee$ $(i,j) \notin \text{next}(\text{Packet_router}(p)) \}$	

LEMMA(5.71) states that a *termination* state is an *acceptable final* state.

St Pk	LEMMA (5.71)
$\text{NETWORK};$ $(\forall (i,j):\text{Network})$ $(\forall p:\text{Pk})$ $(p \notin (\text{Start_in}(i) \cup \text{IN}(i)) - \text{OUT}(i)) \vee$ $(i,j) \notin \text{next}(\text{Packet_router}(p))$ $\vdash \text{FINAL_NW}$	
<p>proof:</p> $(\forall i:\text{St} \mid i \in \text{dom}(\text{Network}))$ $(\forall p:\text{Pk} \mid p \in (\text{Start_in}(i) \cup \text{IN}(i)) - \text{OUT}(i))$ $(\forall j:\text{St} \mid (i,j) \in \text{Network})$ <p>we deduce from the hypothesis</p> $(1) \quad (i,j) \notin \text{next}(\text{Packet_router}(p))$ <p>as p is 'within' the buffer, we conclude</p> $(2) \quad i \in \text{ran}(\text{Packet_router}(p))$ <p>from (1) . (2) we get</p> $(3) \quad i = \text{last}(\text{Packet_router}(p))$ <p>which according to DEF(5.63) completes the proof</p>	

6. DECOMPOSITION.

A complex system is easier to comprehend if it is described through several specifications presenting the system in increasing degrees of detail. For that reason, and in order to make the development process more manageable as well as to make the end-product easier to modify, it has become a convention to describe (if not develop) systems in a *top-down* fashion. A system is first described in terms of a limited set of modules and their inter-relationship. When first presented the modules are specified solely in terms of their *external interfaces*, i.e. the services they provide for the surrounding modules. In a more refined description the decisions made with respect to the *internal structure* of each module will be documented. During the development the designer must verify that the internal structure provided for each module is *consistent* with the definition of its external interface.

The independent module (as described in section 4.2.1) which is specified through an *input-output* relation can be *decomposed* by giving a set of new independent modules. From the description of the input-output relation for each individual sub-module we can infer the input-output relation computed by the compound system, giving us a method for analysing whether the behaviour of the *refined* system agrees with the behaviour of the system it is intended to *implement*. The techniques for calculating the input-output relation computed by a compound system are based on the assumption that each individual module-realisation runs in *isolation*.

As mentioned in section 4.3, it is important to provide a comparable *refinement* and *development* method for systems whose constituents are process-modules. However, these modules, which run in *parallel* and which, during their execution, can *inspect* and *influence* the execution of one another, cannot be described through a simple input-output relation. Hence, the refinements methods and the techniques for validating a decomposition will take a different form from the methods used for modules running in isolation.

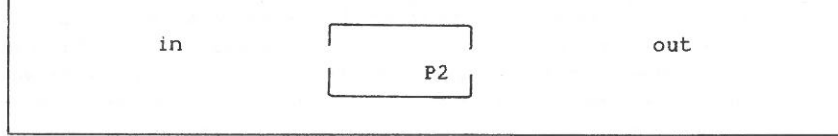
A process-module is refined by replacing one specification of the module with another specification, in which the assertions about the histories of the communication are decomposed and partitioned. A refinement reflects the decision to realise a process-module using two or more concurrently operating process components (a decision taken, for example, in order to increase performance). A refinement will normally increase the number of visible connections as some local communications between the modules of the refinement will be added to the description.

6.1. INTRODUCTORY EXAMPLE.

This section will introduce the decomposition method by means of an example.

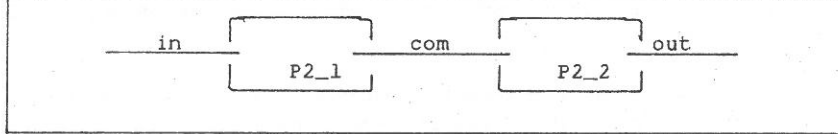
Consider the system presented in section 4.2. (DEF(4.4)) which can be illustrated-

P2_MODULE _____ FIGURE (6.1)



Its decomposed (distributed) counterpart from section 4.3. (DEF (4.9)) can be illustrated

P2_SYSTEM _____ FIGURE (6.2)

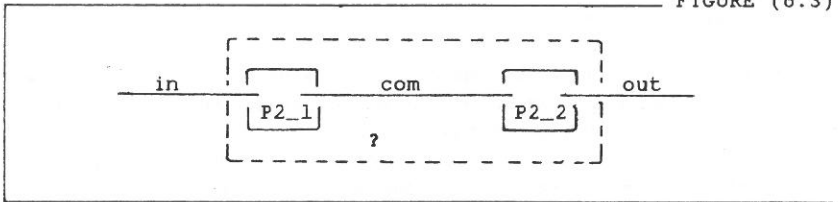


Informally, we will say that a distributed system (like P2_SYSTEM) is an *implementation* (or valid decomposition) of a module (like P2_MODULE) if its externally observable behaviour (i.e. for this example the behaviour along in and out) is indistinguishable from the behaviour of the module it implements. Hence, after decomposing a module, we must compare the behaviour which can be derived from the new compound system (after having *hidden* internal communications) with the expected behaviour of the module. For the example above, we must answer the question, is-

$$? = P2$$

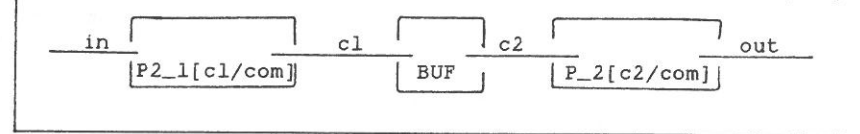
where

_____ FIGURE (6.3)



Before formalising the notion of indistinguishability, consider, for comparison, the following system-

P2_SYSTEM2 _____ FIGURE (6.4)



where

$$BUF[X] \equiv [c1, c2:seq[X] \mid c2 \subseteq c1; \#c1 - \#c2 \leq 1] \quad \text{DEF (6.5)}$$

In an *Environment* where

$$\#in \leq Limit$$

for some *Limit*, and with a final state requirement that

$$out = F \circ in$$

We can prove that all three systems- P2_MODULE, P2_SYSTEM and P2_SYSTEM2 - are *totally correct* systems, i.e. they terminate and are deadlock free. However, the dynamic behaviour of P2_SYSTEM2 is different from the dynamic behaviour of P2_SYSTEM and P2_MODULE, e.g. P2_SYSTEM2 may be observed to consume three consecutive inputs before the activity along the in connection is blocked, while P2_MODULE and P2_SYSTEM can consume at most two before blocking. Hence a system, consisting of P2_SYSTEM (or P2_MODULE) in an environment which insists on performing three inputs before participating in an output, will deadlock.

NB, such an environment could be defined by

$$ENV[X] \equiv [in, out:seq[X] \mid \#in \leq 3 \Rightarrow out=0]$$

If P2_SYSTEM2 is placed in the same environment no deadlock will occur.

As the dynamic behaviour of P2_SYSTEM2 is distinct from that of P2_MODULE, we will *not* accept

$$P2_1[c1/com]; BUF; P2_2[c2/com]$$

as a valid decomposition (or distribution) of

$$P2$$

A method for proving that [P2_1; P2_2] is a valid decomposition of P2 is given in the next section.

6.2. VALID DECOMPOSITION.

Let us consider a system which is specified in terms of predicates on its behaviour along external connections only, i.e.

$$M_SPEC \equiv [\text{ext_con} \mid \text{PRED}_{\text{spec}}] \quad \text{DEF (6.6)}$$

A refinement (decomposition) of this module will expose some internal interfaces and the connections across them. Furthermore a new set of predicates will be given, now involving both the communications along the external connection and the internal connections. We have

$$M_IMP \equiv [\text{ext_con}; \text{int_con} \mid \text{PRED}_{\text{imp}}] \quad \text{DEF (6.7)}$$

Example:

For the system described in the previous section we have a specification of a module (introduced in section 4.3.),

$$P2_SPEC \quad \text{X} \quad \text{DEF (6.8)}$$

$\text{in, out} : \text{seq}[X]$
$\# \text{in} - \# \text{out} \leq 2; \text{out} \subseteq F \circ \text{in}$

and a description of an Implementation of that module

$$P2_IMP \quad \text{X} \quad \text{DEF (6.9)}$$

$\text{in, out, com} : \text{seq}[X]$
$\# \text{in} - \# \text{com} \leq 1; \text{com} \subseteq G \circ \text{in};$
$\# \text{com} - \# \text{out} \leq 1; \text{out} \subseteq H \circ \text{com}$

which is a distributed description because the axioms can be divided in to two groups - one group involving in and com only, and one involving com and out only. Furthermore we claim that this is an *implementation* (valid decomposition) of the specification P2_SPEC.

An analysis of a system consisting of interconnected modules will be carried out without any reference to the internal structure of the modules, i.e. a top-level specification for each module (like M_SPEC) will be used as a basis for the analysis. When refining a system we provide an implementation for each individual module (like M_IMP for M_SPEC) and, in order to ensure that the global dynamic properties (termination and deadlock) are maintained, we must give a guarantee, for each module, that the behaviour of the implementation *simulates* the specified behaviour along the external connections. In other words, the behaviour along the external connections (ext_con) which can be derived from M_IMP must simulate the behaviour along these connections as derived from M_SPEC.

In order to compare the two systems (M_IMP and M_SPEC) we introduce a *projection* function-

$$\text{project} : \{M_IMP\} \rightarrow \{M_SPEC\}$$

NB. this is analogous to a *retrieve-function* [Jones.14] or an *abstraction-function* [Hoare.9]

For the system introduced above we can define the projection function as follows-

$$\text{project} = (\lambda P2_IMP)(\mu P2_SPEC')(\text{in}' = \text{in}; \text{out}' = \text{out}) \quad \text{DEF (6.10)}$$

Let

$$x1, x2 : \text{seq}[Y]$$

be the only external connections for M_SPEC (and M_IMP) then a communication of a value y along the $x1$ connection can be described as-

1) According to the M_SPEC description-

$$\text{comx1}_{\text{spec}}(y) = (\lambda M_SPEC)(\mu M_SPEC')(x1' = x1 * \langle y \rangle; x2' = x2) \quad \text{DEF (6.10)}$$

2) According to the M_IMP description-

$$\text{comx1}_{\text{imp}}(y) = (\lambda M_IMP)(\mu M_IMP')(x1' = x1 * \langle y \rangle; x2' = x2) \quad \text{DEF (6.11)}$$

In the following we will use the examples introduced above to illustrate the method for proving that an implementation is a valid decomposition of its specification.

We must prove-

1) A *Consistency* property.

External activity (as performed by the implementation) must only occur when it is expected to occur, i.e. if a communication (along an external connection) can take place according to the *rules* imposed by the implementation (e.g. M_IMP) then the same communication must also be permitted by the rules imposed by the specification (M_SPEC).

Example 1:

For the communication along $x1$ in the system described by M_IMP, we must prove-

$$(\forall y:Y) (\text{project}(\text{dom}(\text{comx1}_{\text{imp}}(y))) \subseteq \text{dom}(\text{comx1}_{\text{spec}}(y)))$$

which can be verified simply by proving-

$$M_IMP \vdash M_SPEC$$

Intuition - If the axioms of the specification are kept invariant by the operations performed by the implementation then these operations (extension of a history sequence) are also 'allowed' to take place by the specification (NB, an operation can take place whenever its occurrence does not invalidate the axioms).

Example 2:

For the system P2_IMP we must prove-

$$P2_consistency \text{-----} \text{THEOREM (6.13)}$$

$$P2_IMP \vdash \#in - \#out \leq 2; out \subseteq F \cdot in$$

NB:

Because we require that an implementation is *consistent* with its specification we can always describe an implementation as an extension of its specification, e.g. P2_IMP could have been defined as

$$P2_IMP \text{-----} X \ Z \ Y \text{-----} \text{DEF (6.14)}$$

$$P2_SPEC[X, Y];$$

$$com : seq[Z]$$

$$\#in - \#com \leq 1; com \subseteq G \cdot in;$$

$$\#com - \#out \leq 1; out \subseteq H \cdot com$$

2) An *adequacy* property.

When an external communication is expected to take place (according to the rules of the specification), then either this communication or an internal communication can take place (according to the rules of the implementation).

Example 1:

For the communication along x_1 in M_IMP we have-

$$(\forall y:Y)$$

$$(\text{project}^{-1}(\text{dom } com_{x_1}^{spec}(y)) \subseteq$$

$$\text{dom}(com_{x_1}^{imp}(y) \cup \text{dom}(\text{internal_activity})))$$

where *internal_activity* is a relation describing internal communications.

Example 2:

For the system P2_IMP we must prove

$$P2_adequacy \text{-----} \text{THEOREM (6.15)}$$

$$P2_IMP \vdash$$

$$\delta < 2 \Rightarrow \delta_1 < 1 \vee (\delta_1 = 1 \wedge \delta_2 < 1) ;$$

$$\delta \neq 0 \Rightarrow \delta_2 \neq 0 \vee (\delta_2 = 0 \wedge \delta_1 \neq 0)$$

where

$$\delta = \#in - \#out;$$

$$\delta_1 = \#com - \#in;$$

$$\delta_2 = \#out - \#com$$

which states -

a) If an input is expected to occur, i.e.

$$\delta < 2$$

then either the input can be performed, i.e.

$$\delta_1 < 1$$

or an internal communication can take place, i.e.

$$\delta_2 < 1$$

b) If an output is expected to occur, i.e.

$$\delta \neq 0$$

then an output can take place, i.e.

$$\delta_2 \neq 0$$

or an internal communication can take place, i.e.

$$\delta_1 \neq 0$$

3) A *Halting* property.

The adequacy property does not guarantee that a communication which is expected to take place will ever happen, because non-terminating internal activity (*infinite chatter*) could prevent the system from ever reaching a state which allows the occurrence of the external communication.

We must verify that internal activity will terminate if the system is *unstimulated* (no external activities will occur). This property can be verified by presenting a decreasing variant function for the *internal* activity.

Example:

The following function is a variant function for system P2_IMP

$$P2_var = (\lambda P2_IMP)(\#in - \#com) \text{-----} \text{DEF (6.16)}$$

Summary-

The *consistency* property ensures that communications along the external connection only occur when they are expected to occur. The *adequacy* and *halting* properties ensure that a communication which is expected to occur can occur before the

occurrence of any other external communication. Furthermore, the occurrence can only be delayed by a finite number of internal events. (NB, compare these properties with the *readiness* principles from section 4.2.3).

An implementation which satisfies all three properties with respect to a specification is said to be a *valid decomposition* of that specification. Hence, P2_IMP is a *valid decomposition* of P2_SPEC (see THEOREM(6.13), THEOREM(6.15) and DEF(6.16)).

6.3. LIMITATIONS OF THE METHOD.

In this section we will define some systems which consist of two or more process-modules. The examples are given mainly to reveal the limitations of the method for decomposition (or distribution) presented in the two previous sections. Suggestions for surmounting these limitations will be given.

6.3.1. Piped Systems.

As the example (P2) presented in the previous sections illustrated, an *implementation* (or *refinement*) where a single process-module is replaced by a set of process-modules in a *pipe-line* arrangement is easily expressible in the provided framework. Furthermore the tools proved to be adequate for checking the behaviour of such pipe-line arranged systems against an abstract description in which no reference is made to the (local) communications between the individual stages of the pipe-line.

The following example is a continuation of the bounded buffer example from section 5.1.

We will apply the method developed in the previous section to the implementation of a bounded buffer as specified by BB, (DEF(5.3)). The example is similar to the example (P2) presented previously, and does not add to the explanation of the method. However the example is included to give yet another overview of the approach for decomposing systems and to illustrate the convenience of the schema-notation used as a development tool.

let

B1, B2 : N1

be given. In order to verify that

PIPE DEF (6.17)

$$\text{BB}[\text{com}/\text{out}; B1/B]; \text{BB}[\text{com}/\text{in}; B2/B]$$

is an implementation of

$\text{BB}[B1+B2/B]$

we must prove

- 1) a *consistency* theorem, which is

THEOREM (6.18)

$$\text{PIPE} \vdash \text{BB}[B1+B2/B]$$

- 2) an *adequacy* theorem, which for this system can take the form

THEOREM (6.19)

```

PIPE;
#in - #com = 0  $\vee$  #com - #out = B2
 $\vdash$ 
 $\delta < B1+B2 \Rightarrow \delta 1 < B1;$ 
 $\delta \neq 0 \Rightarrow \delta 2 \neq 0$ 
where
 $\delta = \#in - \#out;$ 
 $\delta 1 = \#in - \#com;$ 
 $\delta 2 = \#com - \#out$ 

```

- 3) a *halting* property, which is ensured by the presentation of a variant function,

$\text{Var}_{\text{Pipe}} = (\lambda \text{ PIPE})(\#in - \#com)$

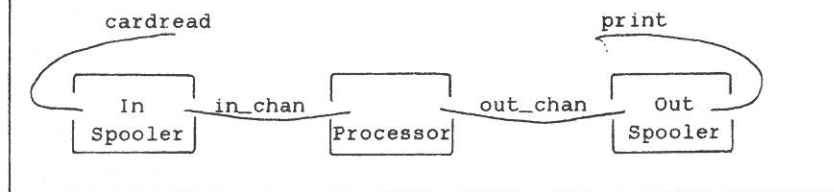
6.3.2. Unordered Pipe.

The unordered queue introduced in section 5.2.2. (DEF(5.17)) can, as mentioned in that section, conveniently represent queues of jobs in an operating system.

let

In_Spooler[J] \equiv QUEUE[J][cardread/arrive; in_chan/depart];
 Processor[J] \equiv QUEUE[J][in_chan/arrive; out_chan/depart];
 Out_Spooler[J] \equiv QUEUE[J][out_chan/arrive; print/depart]

be components of a simplified batch operating system, which can be represented as

BATCH FIGURE (5.20)

BATCH is indeed an implementation of SYSTEM, where

$\text{SYSTEM} \equiv \text{QUEUE}[\text{cardread}/\text{arrive}; \text{print}/\text{depart}]$ DEF (5.21)

let us now consider a slightly modified version of BATCH. By restricting the buffering capacity of the Out_Spooler, i.e. for some $L : \mathbb{N}$ we may have

$\text{Out_Spooler2} \equiv [\text{Out_Spooler} \mid \# \text{out_chan} - \# \text{print} \leq L]$

and by defining a new batch system as

$\text{BATCH2} \equiv [\text{In_Spooler}; \text{Processor}; \text{Out_Spooler2}]$ DEF (5.22)

we get a system whose behaviour *cannot* be explained by an abstraction which does not refer to the internal communications between Processor and Out_Spooler. BATCH2 is indeed distinct from BATCH because BATCH allows for any job currently within the system to get priority and to run to completion (i.e. being served by the printer), while this may not be possible for system BATCH2. In a state where

$$\# \text{out_chan} - \# \text{print} = L$$

BATCH2 cannot accept further requests along out_chan before another job has left the print-spooler queue. Hence a job currently in the input-spooler queue or in the processor queue cannot be completed before at least one other job has been completed.

BATCH2 therefore (according to the rules stated in the previous section) is not an implementation of SYSTEM. However if we can *rely* on the environment never to insist on the completion of a particular job, then, intuitively, BATCH2 will be a perfect implementation of SYSTEM. This can be formalised.

Let $\text{O_ENV}(p)$ describe the environment's restriction on the printer connection (p). If we can prove

THEOREM (5.23)

$$\begin{aligned}
 & (\forall i:N; j, j':J) \\
 & (\forall p_s : \text{Print_Stream}) \\
 & ((i \rightarrow j) \in p_s) \Rightarrow ((p_s \oplus \{i \rightarrow j'\}) \in \text{Print_Stream}) \\
 & \text{where} \\
 & \text{Print_Stream} = \{p : \text{seq}[J] \mid \text{O_ENV}(p)\}
 \end{aligned}$$

which states that the given environment does not insist on the completion of any particular job, then we will accept BATCH2 as an implementation of SYSTEM.

NB, if we do not make this (or a similar) assumption about the environment then BATCH2 cannot be considered to be a refinement of SYSTEM because BATCH2 might reject requests which could be accepted by SYSTEM, i.e. introducing BATCH2 as a refinement for SYSTEM may also introduce deadlock possibilities.

6.3.3. A Resource Monitor.

Consider the system in DEF (5.44) whose constituents are R_1 from DEF (5.37), R_2 from DEF (5.38) and ENV from DEF (5.45).

The subsystem

$\text{RL_PIPE} \equiv [R_1; R_2]$ DEF (6.24)

does not implement

$\text{RL_SYS} \xrightarrow{\quad R \ L \quad} \text{DEF (6.25)}$

$$\begin{aligned}
 & \text{in}_R, \text{out}_R : \text{seq}[R]; \\
 & \text{in}_L, \text{out}_L : \text{seq}[L] \\
 & \text{---} \\
 & \text{out}_R \subseteq \text{in}_R; \text{out}_L \subseteq \text{in}_L; \\
 & \delta_R + \delta_L \leq 2 \\
 & \text{where} \\
 & \delta_R = \# \text{in}_R - \# \text{out}_R; \\
 & \delta_L = \# \text{in}_L - \# \text{out}_L
 \end{aligned}$$

as the RL_PIPE will deadlock in the situation where

$$\delta_1 = 1 \wedge \delta_2 = 1$$

However, the system

RL_PIPE2 _____ DEF (6.26)

RL_PIPE

$$\delta_R = 0 \vee \delta_L = 0$$

where

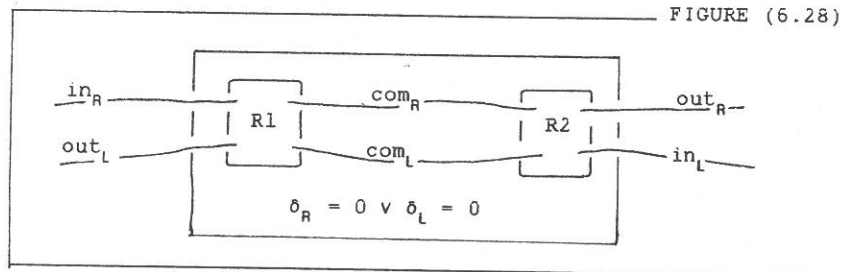
$$\delta_R = \#in_R - \#out_R;$$

$$\delta_L = \#in_L - \#out_L$$

does implement

RL_SYS2 \equiv [RL_SYS | $\delta_R=0 \vee \delta_L=0$] DEF (6.27)

RL_PIPE2 is a system consisting of three processes, one of which (a Monitor) does not refer to the local communications between the resource controllers R1 and R2. This system can be illustrated as



7. CONCLUSION.

The goal of this thesis has been to develop a method for *specifying, analysing and refining* the designs of distributed systems.

A new approach for *specifying* distributed systems has been proposed in chapter 4 to accomplish this goal. The approach can be used to specify distributed systems in different degrees of detail; i.e. a description with few *visible internal interfaces* is considered to be more abstract than a description in which all interfaces from the future realisation take part. The systems are described through *predicates* on the *histories* of communications across the visible interfaces, and the modularisation (or distribution) of the system is determined by the possible grouping of these predicates (*localisation of knowledge*).

A method for *analysing* the behaviour of distributed designs has been proposed in sections 4.5 and 5.2. The analysis is based on a technique presented in chapter 3. It should be pointed out that the formal framework developed (in chapter 3) for analysing *non-deterministic* systems consists mainly of reformulations of familiar concepts (transition-systems: [Keller,16], [Lamswerde,17]; well-founded relations: [Abrial,1]; and correctness of non-deterministic constructs: [Dijkstra,7]). The techniques presented in chapter 3 are, as illustrated in section 3.2. (A Bounded Buffer) applicable to descriptions based on abstractions of internal states as well as to descriptions based on abstractions of external behaviour (section 5.1. A Bounded Buffer).

An important aspect of the method presented is its use in *refinement* or *step-wise development* of distributed systems. Chapter 6 introduced a simple method for comparing different levels of descriptions using projections as *retrieve-functions* [Jones,14]. It was demonstrated that the refinement method is especially suitable for refinements where a single process-module is replaced by several modules arranged in a *pipe-line*. Note that when abstracting details of the protocols (*hand-shaking, synchronisations* etc.) used for communications between process-components, the number of systems falling into this category is quite large. Some limitations of the method have been pointed out, together with informal suggestions, for overcoming them. From example 6.3.2.(Unordered Pipe) we learned that some *internal interfaces* (i.e. internal communications) of compound systems can be so integral to the system's behaviour that they cannot be *hidden* without losing information about the true behaviour along the external connections of the system. Instead of hiding such internal events (and introducing an unmanageable kind of non-determinism) a 'solution', in which assumptions are made about the behaviour of future environments, is proposed. (*rely-conditions*).

7.1. RELATED WORK.

Communication between cooperating processes has become an area of concern especially since the introduction of multiprogramming systems. In such systems the possibility for breaking down large tasks into smaller communicating (concurrently operating) processes provides a useful abstraction tool. The advance in technology - i.e. the introduction of inexpensive microprocessors - has further stimulated research in this area.

The basis for a formal treatment of cooperating processes was laid with the introduction of *semaphores* [Dijkstra,6]. Since then, several programming languages which embody the notion of processes have been suggested and high level language constructs, which provide the programmer with an abstract (and 'secure') tool for managing process communication and synchronisation, have been proposed. Examples of earlier work in this area are to be found in the work of [Hoare,10] and [Brinch Hansen,2] on *monitors* and [Campbell,3] on *path expressions*. More recently and certainly more notably, there are the high-level communication facilities provided through CSP [Hoare,11] and Ada [13].

Verification of concurrent programs has also been investigated (e.g. [Owicki,19]). Verification of programs requires a *specification* with which the program can be compared, which suggests that it is important to have a framework (or a style) in which specifications can be expressed. The methods for specifying *sequential* programs are well understood; however, methods for specifying *distributed* (concurrent) programs are still in their infancy. For distributed systems which consist of several cooperating components the properties to be specified consist not only in WHAT is to be computed but also WHEN and WHERE (see section 4.3). Hence an approach for specifying - like the approach presented in this thesis - *must* be concerned with *time* in one form or another.

In this thesis sequences of past communications (i.e. histories of communications) across interfaces are used to give the needed notion of *time*.

In the area of *program verification* we find the use of *communication sequences* as *ghost* (or mythical) variables introduced in order to verify the correctness of *co-routines* [Clint,5]. Assertions involving the past history of communications along external channels have more recently been used in [Chen,4] and [Hoare,12] (NB, the assertions used in [Hoare,12] also involve the immediate future).

An alternative approach for documenting *time-related* properties for concurrent systems is to be found in the work on *temporal logic* (e.g. [Manna,18]). In this approach the time-related properties are classified into: *invariant* properties, *eventuality* properties and *precedence* properties; these properties can all be expressed using predicates on communication sequences.

The main goal of the work presented in this thesis is *not* to provide a method for *program proving*, but rather to provide a *development* tool for *distributed systems*, i.e. to provide a uniform formal notation in which design decisions can be

documented and analysed, and a method for refining and modifying such systems. Little work has been done in this area. However, notable suggestions have recently been made in [Jones,15] where a 'module' executing in an interfering environment is described using *guarantee*- and *rely* assertions on the information it *shares* with other modules.

7.2. FUTURE RESEARCH.

Because the approach presented in this thesis has only been under development for a short time, research into feasibility of the proposed methods cannot be declared complete, nor can the methods themselves.

The tools must be studied further and I am convinced that inspiration for further research can be obtained from 1) [Hoare,12] in which assertions on *past* communications are used to describe intended behaviour for processes, 2) [Jones,15] where the goal is to find development-methods for concurrent systems, and 3) [Manna,18] who uses *temporal logic* in order to abstract direct references to execution sequences.

One aspect of the development process which has been intentionally omitted from this thesis is the problem of *realisation*. The difficulties in realising a communication between two distinct process components have not been considered. In any development process we must at some level of refinement consider the physical properties of the target machinery. Hence, for distributed systems which consist of communicating process-components we must consider the means of communication provided by the inter-process control mechanisms built into the realisation tool.

The abstract communications of the operational model in terms of which we analyse our specifications may not be easily realisable. Intuitively the difficulties will depend on the realisation tool (compare the facilities provided by Ada or CSP with communications performed entirely by means of shared memory and semaphore-operations). We may reformulate a specification before realisation in order to simplify the task of translating the specification into a program.

Example:

Consider the communications between Q and CS in FIGURE(5.28). At any point in time Q will be willing to communicate any 'job' which does not invalidate the axioms of Q, e.g.

$$\{ j_1, j_2, j_3, j_4 \}$$

and at any point in time CS is willing to 'receive' the set of jobs which will not invalidate the *mutual* exclusion property, e.g.

$$\{ j_2, j_3 \}$$

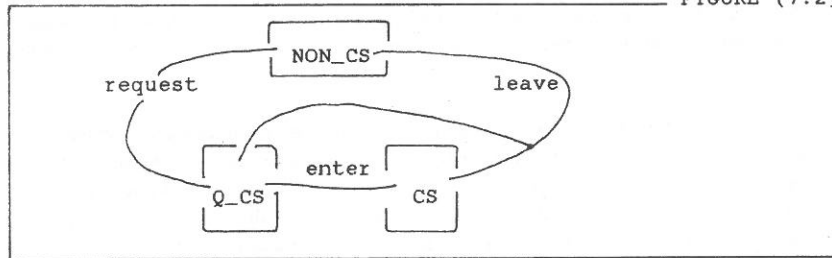
If the intersection of the two sets is non-empty then a communication will take place. Such a communication is not directly realisable. However, we can *reformulate* the specification and obtain a description which is 'closer' to a realisation.

Let us define (use DEF(5.32) for Q)

Q_CS	J	DEF(7.1.)
Q ; leave : seq[J]		
$M(\text{inside}) \cap \text{inside} = \{\}$ where inside = { $j:J \mid \text{card}(\text{enter}\downarrow\{j\}) > \text{card}(\text{leave}\downarrow\{j\})$ }		

as replacement for Q , which gives

FIGURE (7.2)



NB. the modification does *not* change the properties of the system as no axioms have been added. Furthermore the mutual exclusion axiom from CS can now be removed as this axiom will be 'obeyed' by Q_CS . Hence, the communications between Q_CS and CS are much simpler as all knowledge about the communications along *enter* has been localised within Q_CS .

Further investigations must be made in the area of realisation of specifications.

It is most pressing, however, to apply the methods presented to a wider range of distributed systems, not only in order to improve the methods but also to explore their limitations and to identify the category of systems for which they are suitable.

List of References.

- [1] J.R. Abrial:
'The Specification Language Z: Basic Library'.
Report from Oxford University Computing Laboratory, 1980.
- [2] P. Brinch Hansen:
'Operating System Principles'
Prentice-Hall, 1973.
- [3] R.H. Campbell and A.N. Habermann:
'The Specification of Process Synchronization by Path Expressions'.
Lecture notes in Comp.Sci. Vol 16.
Springer Verlag, 1974.
- [4] Zhou Chao Chen and C.A.R. Hoare:
'Partial Correctness of Communicating Processes and Protocols'.
Proc. International Conference on Distributed Computing, 1981.
- [5] M. Clint:
'Program proving: Coroutines'.
Acta Informatica 2, 1973.
- [6] E.W. Dijkstra:
'Cooperating Sequential Processes'.
Programming Languages, Academic Press, New York, 1968.
- [7] E.W. Dijkstra:
'Guarded Commands, Non-determinacy and Formal Derivation of Programs'.
CACM Vol 18 No 8, 1975.
- [8] E.W. Dijkstra:
'A Discipline of Programming'.
Prentice-Hall Int. Series in Automatic Computation, 1976.
- [9] C.A.R. Hoare:
'Proof of Correctness of Data Representations'.
Acta Informatica 1, 1972.
- [10] C.A.R. Hoare:
'Monitors: An Operating System Structuring Concept'.
CACM Vol 17 No 10, 1974.
- [11] C.A.R. Hoare:
'Communicating Sequential Processes'.
CACM Vol 21 No 8, 1978.
- [12] C.A.R. Hoare:
'A Calculus of Total Correctness for Communicating processes'.
Oxford University Computing Laboratory, PRG Monograph no 23, 1981.

- [13] Reference Manual.
'Rationale for the Design of the GREEN Programming Language'
Honeywell, inc. Minneapolis. 1979.
- [14] C.B. Jones:
'Software Developmant - A Rigorous Approach'.
Prentice-Hall Int. Series in Computer Science, 1980.
- [15] C.B. Jones:
'Development method of Computer Programs Including a Notlon of
Interference'
D.Phil. thesis, University of Oxford. 1981.
- [16] R.M. Keller:
'Formal Verification of Parallel Programs'.
CACM Vol 19 No 7. 1976.
- [17] A. von Lamswerde and M. Sintzoff:
'Formal Derivation of Strongly Correct Concurrent Programs'.
Acta Informatica 12. 1979.
- [18] Z.Manna and A. Pnuelli:
'Verification of Concurrent Programs: The Temporal Framework'.
Lecture notes for The International Summer School, Munich, August
1981.
- [19] S. Owicki and D Gries:
'Verifying properties of Parallel Programs: An Axiomatic Approach'.
CACM Vol 19 No 5. 1976.
- [20] P.L. Parnas:
'On the Criteria to be used in Decomposing Systems into Modules'.
CACM Vol 15 No 12. 1972.
- [21] B. Sufrin:
'Formal Specification of a Display Editor'.
Oxford University Computing Laboratory. PRG Monograph no 21. 1981.