

# ABSTRACT INTERPRETATION : A UNIFIED LATTICE MODEL FOR STATIC ANALYSIS OF PROGRAMS BY CONSTRUCTION OR APPROXIMATION OF FIXPOINTS

Patrick Cousot\* and Radhia Cousot\*\*

Laboratoire d'Informatique, U.S.M.G., BP. 53  
38041 Grenoble cedex, France

## 1. Introduction

A program denotes computations in some universe of objects. Abstract interpretation of programs consists in using that denotation to describe computations in another universe of abstract objects, so that the results of abstract execution give some informations on the actual computations. An intuitive example (which we borrow from Sintzoff [72]) is the rule of signs. The text  $-1515 * 17$  may be understood to denote computations on the abstract universe  $\{(+), (-), (\pm)\}$  where the semantics of arithmetic operators is defined by the rule of signs. The abstract execution  $-1515 * 17 \Rightarrow -(+) * (+) \Rightarrow (-) * (+) \Rightarrow (-)$ , proves that  $-1515 * 17$  is a negative number. Abstract interpretation is concerned by a particular underlying structure of the usual universe of computations (the sign, in our example). It gives a summary of some facets of the actual executions of a program. In general this summary is simple to obtain but inaccurate (e.g.  $-1515 + 17 \Rightarrow -(+) + (+) \Rightarrow (-) + (+) \Rightarrow (\pm)$ ). Despite its fundamentally incomplete results abstract interpretation allows the programmer or the compiler to answer questions which do not need full knowledge of program executions or which tolerate an imprecise answer, (e.g. partial correctness proofs of programs ignoring the termination problems, type checking, program optimizations which are not carried in the absence of certainty about their feasibility, ...).

## 2. Summary

Section 3 describes the syntax and mathematical semantics of a simple flowchart language, Scott and Strachey[71]. This mathematical semantics is used in section 4 to build a more abstract model of the semantics of programs, in that it ignores the sequencing of control flow. This model is taken to be the most concrete of the abstract interpretations of programs. Section 5 gives the formal definition of the abstract interpretations of a program.

\* Attaché de Recherche au C.N.R.S., Laboratoire Associé n° 7.

\*\* This work was supported by IRIA-SESORI under grants 75-035 and 76-160.

Abstract program properties are modeled by a complete semilattice, Birkhoff[61]. Elementary program constructs are locally interpreted by order preserving functions which are used to associate a system of recursive equations with a program. The program global properties are then defined as one of the extreme fixpoints of that system, Tarski[55]. The abstraction process is defined in section 6. It is shown that the program properties obtained by an abstract interpretation of a program are consistent with those obtained by a more refined interpretation of that program. In particular, an abstract interpretation may be shown to be consistent with the formal semantics of the language. Levels of abstraction are formalized by showing that consistent abstract interpretations form a lattice (section 7). Section 8 gives a constructive definition of abstract properties of programs based on constructive definitions of fixpoints. It shows that various classical algorithms such as Kildall [73], Wegbreit[75] compute program properties as limits of finite Kleene[52]'s sequences. Section 9 introduces finite fixpoint approximation methods to be used when Kleene's sequences are infinite, Cousot[76]. They are shown to be consistent with the abstraction process. Practical examples illustrate the various sections. The conclusion points out that abstract interpretation of programs is a unified approach to apparently unrelated program analysis techniques.

## 3. Syntax and Semantics of Programs

We will use finite flowcharts as a language independent representation of programs.

### 3.1 Syntax of a Program

A program is built from a set "Nodes". Each node has successor and predecessor nodes :

$$\begin{aligned} \underline{n\text{-succ}}, \underline{n\text{-pred}} : \text{Nodes} \rightarrow 2^{\text{Nodes}} \mid (n \in \underline{n\text{-succ}}(n)) \\ \Leftrightarrow (n \in \underline{n\text{-pred}}(m)) \end{aligned}$$

Hereafter, we note  $|S|$  the cardinality of a set  $S$ . When  $|S| = 1$  so that  $S = \{x\}$  we sometimes use  $S$  to denote  $x$ .

The node subsets "Entries", "Assignments", "Tests", "Junctions" and "Exits" partition the set Nodes.

- An entry node ( $n \in \text{Entries}$ ) has no predecessors and one successor, ( $(\underline{n\text{-pred}}(n) = \emptyset)$  and  $(|\underline{n\text{-succ}}(n)| = 1)$ ).

- An assignment node ( $n \in \text{Assignments}$ ) has one predecessor and one successor ( $(|\underline{n\text{-pred}}(n)| = 1)$  and  $(|\underline{n\text{-succ}}(n)| = 1)$ ). Let "Ident" and "Expr" be the distinct syntactic categories of identifiers and expressions. An assignment node  $n$  assigns the value of the right hand-side expression  $\underline{\text{expr}}(n)$  to the left hand-side identifier  $\underline{\text{id}}(n)$ :  

$$\underline{\text{expr}} : \text{Assignments} \rightarrow \text{Expr}$$

$$\underline{\text{id}} : \text{Assignments} \rightarrow \text{Ident}$$
- A test node ( $n \in \text{Tests}$ ) has a predecessor and two successors, ( $(|\underline{n\text{-pred}}(n)| = 1)$  and  $(|\underline{n\text{-succ}}(n)| = 2)$ ). The true and false successor nodes are respectively denoted  $\underline{n\text{-succ-t}}(n)$  and  $\underline{n\text{-succ-f}}(n)$ :  

$$\underline{n\text{-succ-t}}, \underline{n\text{-succ-f}} : \text{Tests} \rightarrow \text{Nodes}$$

$$(\forall n \in \text{Tests}, \underline{n\text{-succ}}(n) = \{\underline{n\text{-succ-t}}(n), \underline{n\text{-succ-f}}(n)\}).$$

Let "Bexpr" be the syntactic category of boolean expressions, each test node  $n$  contains a boolean expression  $\underline{\text{test}}(n)$ :  

$$\underline{\text{test}} : \text{Tests} \rightarrow \text{Bexpr}$$
- A junction node ( $n \in \text{Junctions}$ ) has one successor and more than one predecessor, ( $(|\underline{n\text{-succ}}(n)| = 1)$  and  $(|\underline{n\text{-pred}}(n)| > 1)$ ). Immediate predecessor nodes of a junction node are not junction nodes, ( $\forall n \in \text{Junctions}, \forall m \in \underline{n\text{-pred}}(n), \text{not}(m \in \text{Junctions})$ ).
- An exit node  $n$  has one predecessor and no successor, ( $(|\underline{n\text{-pred}}(n)| = 1)$  and  $(\underline{n\text{-succ}}(n) = \emptyset)$ ).

The set "Arcs" of edges of a program is a subset of  $\text{Nodes} \times \text{Nodes}$  defined by:

$$\text{Arcs} = \{\langle n, m \rangle \mid (n \in \text{Nodes}) \text{ and } (m \in \underline{n\text{-succ}}(n))\}$$

which may be equivalently defined by:

$$\text{Arcs} = \{\langle n, m \rangle \mid (m \in \text{Nodes}) \text{ and } (n \in \underline{n\text{-pred}}(m))\}.$$

We will assume that the directed graph  $\langle \text{Nodes}, \text{Arcs} \rangle$  is connected.

We will use the following functions:

$$\underline{\text{origin}}, \underline{\text{end}} : \text{Arcs} \rightarrow \text{Nodes} \mid (\forall a \in \text{Arcs}, a = \langle \underline{\text{origin}}(a), \underline{\text{end}}(a) \rangle)$$

$$\underline{a\text{-succ}} : \text{Nodes} \rightarrow 2^{\text{Arcs}} \mid$$

$$\underline{a\text{-succ}}(n) = \{\langle n, m \rangle \mid m \in \underline{n\text{-succ}}(n)\}$$

$$\underline{a\text{-pred}} : \text{Nodes} \rightarrow 2^{\text{Arcs}} \mid$$

$$\underline{a\text{-pred}}(n) = \{\langle m, n \rangle \mid m \in \underline{n\text{-pred}}(n)\}$$

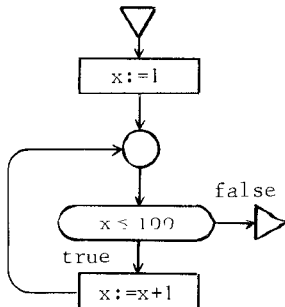
$$\underline{a\text{-succ-t}} : \text{Tests} \rightarrow \text{Arcs}$$

$$\underline{a\text{-succ-t}}(n) = \langle n, \underline{n\text{-succ-t}}(n) \rangle$$

$$\underline{a\text{-succ-f}} : \text{Tests} \rightarrow \text{Arcs}$$

$$\underline{a\text{-succ-f}}(n) = \langle n, \underline{n\text{-succ-f}}(n) \rangle$$

Example:



### 3.2 Semantics of Programs

This section develops a simple "mathematical semantics" of programs, in the style of Scott and Strachey[71].

- If  $S$  is a set we denote  $S^0$  the complete lattice obtained from  $S$  by adjoining  $\{ \perp_S, \top_S \}$  to it, and imposing the ordering  $\perp_S \leq x \leq \top_S$  for all  $x \in S$ .
- The semantic domain "Values" is a complete lattice which is the sum of the lattice  $\text{Bool} = \{\text{true}, \text{false}\}^0$  and some other primitive domains.
- Environments are used to hold the bindings of identifiers to their values:  

$$\text{Env} = \text{Ident}^0 \rightarrow \text{Values}$$
- We assume that the meaning of an expression  $\underline{\text{expr}} \in \text{Expr}$  in the environment  $e \in \text{Env}$  is given by  $\underline{\text{val}} \llbracket \underline{\text{expr}} \rrbracket (e)$  so that:  

$$\underline{\text{val}} : \text{Expr} \rightarrow [\text{Env} \rightarrow \text{Values}].$$

In particular the projection  $\underline{\text{val}} \mid \text{Bexpr}$  of the function  $\underline{\text{val}}$  in domain  $\text{Bexpr}$  has the functionality:  

$$\underline{\text{val}} \mid \text{Bexpr} : \text{Bexpr} \rightarrow [\text{Env} \rightarrow \text{Bool}].$$
- The state set "States" consists of the set of all information configurations that can occur during computations:  

$$\text{States} = \text{Arcs}^0 \times \text{Env}.$$

A state ( $s \in \text{States}$ ) consists in a control state ( $\underline{\text{cs}}(s)$ ) and an environment ( $\underline{\text{env}}(s)$ ), such that:  

$$\forall s \in \text{States}, s = \langle \underline{\text{cs}}(s), \underline{\text{env}}(s) \rangle.$$

- We use a continuous conditional function  $\underline{\text{cond}}(b, e_1, e_2)$  equal to  $\perp$ ,  $e_1$ ,  $e_2$  or  $\top$  respectively as the value of  $b$  is  $\perp$ ,  $\text{true}$ ,  $\text{false}$  or  $\top$ . We also use  $\underline{\text{if}} \ b \ \underline{\text{then}} \ e_1 \ \underline{\text{else}} \ e_2 \ \underline{\text{fi}}$  to denote  $\underline{\text{cond}}(b, e_1, e_2)$ .

- If  $e \in \text{Env}$ ,  $v \in \text{Values}$ ,  $x \in \text{Ident}$  then  

$$\underline{e} [v/x] = \lambda y. \underline{\text{cond}}(y = x, v, \underline{e}(y)).$$

- The state transition function defines for each state  $a$  next state (we consider deterministic programs):  

$$\underline{n\text{-state}} : \text{States} \rightarrow \text{States}$$

$$\underline{n\text{-state}}(s) =$$

$$\underline{\text{let}} \ n \ \underline{\text{be}} \ \underline{\text{end}}(\underline{\text{cs}}(s)), \ e \ \underline{\text{be}} \ \underline{\text{env}}(s) \ \underline{\text{within}}$$

$$\underline{\text{case}} \ n \ \underline{\text{in}}$$

$$\text{Assignments} \Rightarrow \langle \underline{a\text{-succ}}(n), e[\underline{\text{val}} \llbracket \underline{\text{expr}}(n) \rrbracket (e) / \underline{\text{id}}(n)] \rangle$$

$$\text{Tests} \Rightarrow \langle \underline{a\text{-succ-t}}(n), e, \underline{a\text{-succ-f}}(n), e \rangle$$

$$\text{Junctions} \Rightarrow \langle \underline{a\text{-succ}}(n), e \rangle$$

$$\text{Exits} \Rightarrow s$$

$$\underline{\text{esac}}$$

(Each partial function  $f$  on a set  $S$  is extended to a continuous total function on the corresponding domain  $S^0$  by  $f(\perp) = \perp$ ,  $f(\top) = \top$  and  $f(x) = \perp$  if the partial function is undefined at  $x$ ).

- Let  $\perp_{\text{Env}}$  be the bottom function on  $\text{Env}$  such that  

$$(\forall x \in \text{Ident}^0, \perp_{\text{Env}}(x) = \perp_{\text{Values}}).$$

Let  $I$ -states be the subset of initial states:  

$$I\text{-states} = \{ \langle \underline{a\text{-succ}}(m), \perp_{\text{Env}} \rangle \mid m \in \text{Entries} \}$$

- A "computation sequence" with initial state  $i_s \in I\text{-states}$  is the sequence :  

$$s_n = \text{n-state}^n(i_s) \quad \text{for } n = 0, 1, \dots$$
where  $f^0$  is the identity function and  $f^{n+1} = f \circ f^n$ .
- The initial to final state transition function :

$$\text{n-state}^\infty : \text{States} \rightarrow \text{States}$$

is the minimal fixpoint of the functional :  

$$\lambda F. (\text{n-state} \circ F)$$

Therefore

$$\text{n-state}^\infty = Y_{\text{States} \rightarrow \text{States}}(\lambda F. (\text{n-state} \circ F))$$

where  $Y_D(f)$  denotes the least fixpoint of  
 $f : D \rightarrow D$ , Tarski[55].

#### 4. Static Semantics of Programs

The constructive or operational semantics of programs defined in section 3 considers the *sequence* in which states occur during execution. The fundamental remark of Floyd[67] is that to prove static properties of programs it is often sufficient to consider the *sets* of states associated with each program point.

Hence, we define the context  $C_q$  at some program point  $q \in \text{Arcs}$  of a program  $P$  to be the set of all environments which may be associated to  $q$  in all the possible computation sequences of  $P$  :

$$C_q \in \text{Contexts} = 2^{\text{Env}}$$

$$C_q = \{e \mid (\exists n \geq 0, \exists i_s \in I\text{-states} \mid \langle q, e \rangle = \text{n-state}^n(i_s))\}$$

The context vector  $C_v$  associates a context to each of the program points of a program :

$$C_v \in \text{Context-Vectors} = \text{Arcs}^0 \rightarrow \text{Contexts}$$

$$C_v = \lambda q. \{e \mid (\exists n \geq 0, \exists i_s \in I\text{-states} \mid \langle q, e \rangle = \text{n-state}^n(i_s))\}$$

According to the semantics of programs, the context  $C_v(r)$  associated to arc  $r$  is related to the contexts  $C_v(q)$  at arcs  $q$  adjacent to  $r$ , ( $\text{end}(q) = \text{origin}(r)$ ,  $q \rightarrow r$ ). From the definition of the state transition function we can prove the equation :

$$C_v(r) = \text{n-context}(r, C_v)$$

where

$$\text{n-context} : \text{Arcs}^0 \times \text{Context-Vectors} \rightarrow \text{Contexts}$$

is defined by :

$$\text{n-context}(r, C_v) = \text{case } \text{origin}(r) \text{ in}$$

$$\text{Entries} \Rightarrow \{ \perp_{\text{Env}} \}$$

$$\text{Assignments} \cup \text{Tests} \cup \text{Junctions} \Rightarrow$$

$$\bigcup_{q \in \text{a-pred}(\text{origin}(r))} \bigcup_{e \in C_v(q)} \text{env-on}(r)(\text{n-state}(\langle q, e \rangle))$$

$$\text{esac}$$

(We define  $\text{env-on} : \text{Arcs}^0 \rightarrow [\text{States} \rightarrow 2^{\text{Env}}]$  to be  
 $\lambda r. (\lambda s. \text{cond}(r = \text{cs}(s), \{\text{env}(s)\}, \emptyset))$ ).

Since the equation  $C_v(r) = \text{n-context}(r, C_v)$  must be valid for each arc,  $C_v$  is a solution to the system of "forward" equations :

$$C_v = \text{F-cont}(C_v)$$

where

$$\text{F-cont} : \text{Context-Vectors} \rightarrow \text{Context-Vectors}$$

is defined by :

$$\text{F-cont}(C_v) = \lambda r. \text{n-context}(r, C_v)$$

$\text{Context-Vectors}$  is a complete lattice with union  $\tilde{\cup}$  such that  $C_{v_1} \tilde{\cup} C_{v_2} = \lambda r. (C_{v_1}(r) \cup C_{v_2}(r))$ .

$\text{F-cont}$  is order preserving for the ordering  $\tilde{\subseteq}$  of  $\text{Context-Vectors}$  which is defined by :

$$\{C_{v_1} \tilde{\subseteq} C_{v_2}\} \iff \{\forall r \in \text{Arcs}, C_{v_1}(r) \subseteq C_{v_2}(r)\}$$

Hence it is known that  $\text{F-cont}$  has fixpoints, Tarski[55]. However, it is trivial to exhibit examples which show that these fixpoints are not always unique. Fortunately, it can be shown that  $C_v$  is included in any solution  $S$  to the system of equations  $X = \text{F-cont}(X)$ , ( $C_v \tilde{\subseteq} S$ ). Tarski[55] shows that this property uniquely determines  $C_v$  as the least fixpoint of  $\text{F-cont}$ . Thus  $C_v$  can be equivalently defined by :

$$D1 : C_v = \lambda q. \{e \mid (\exists n \geq 0, \exists i_s \in I\text{-states} \mid \langle q, e \rangle = \text{n-state}^n(i_s))\}$$

or

$$D2 : C_v = Y_{\text{Context-Vectors}}(\text{F-cont})$$

The concrete context vector  $C_v$  is such that for any program point  $q \in \text{Arcs}$  of the program  $P$ ,

( $\alpha$ )  $C_v(q)$  contains at least the environments  $e$  which may be associated to  $q$  during any execution of  $P$  :

$$\{\exists i \geq 0, \exists i_s \in I\text{-states} \mid \langle q, e \rangle = \text{n-state}^i(i_s)\} \implies \{e \in C_v(q)\}$$

( $\beta$ )  $C_v(q)$  contains only the environments  $e$  which may be associated to  $q$  during an execution of  $P$  :

$$\{e \in C_v(q)\} \implies \{\exists i \geq 0, \exists i_s \in I\text{-states} \mid \langle q, e \rangle = \text{n-state}^i(i_s)\}$$

$C_v$  is merely a static summary of the possible executions of the program. However, our definitions D1 or D2 of  $C_v$  cannot be utilized at compile time since the computation of  $C_v$  consists in fact in running the program (for all the possible input data). In practice compilers may consider states which can never occur during program execution (e.g. some compilers consider that any program may always perform a division by zero although this is not the case for most programs). Hence compilers may use "abstract" contexts satisfying ( $\alpha$ ) but not necessarily ( $\beta$ ), which therefore correctly approximate the concrete contexts we considered until now.

#### 5. Abstract Interpretation of Programs

##### 5.1 Formal Definition

An abstract interpretation  $I$  of a program  $P$  is a tuple

$$I = \langle A\text{-Cont}, \circ, \leq, \tau, \perp, \text{Int} \rangle$$

where the set of abstract contexts is a complete o-semilattice with ordering  $\leq$ , ( $\{x \leq y\} \iff \{x \circ y = y\}$ ). This implies that  $A\text{-Cont}$  has a supremum  $\tau$ . We suppose also  $A\text{-Cont}$  to have an infimum  $\perp$ .

This implies that A-Cont is in fact a complete lattice, but we need only one of the two join and meet operations. The set of context vectors is defined by  $\widetilde{\text{A-Cont}} = \text{Arcs}^0 \rightarrow \text{A-Cont}$ .

Whatever  $(\text{Cv}', \text{Cv}'') \in \text{A-Cont}^2$  may be, we define :

$$\text{Cv}' \circ \text{Cv}'' = \lambda r. \text{Cv}'(r) \circ \text{Cv}''(r)$$

$$\text{Cv}' \leq \text{Cv}'' = \{\forall r \in \text{Arcs}^0, \text{Cv}'(r) \leq \text{Cv}''(r)\}$$

$$\widetilde{\tau} = \lambda r. \tau \text{ and } \widetilde{\perp} = \lambda r. \perp$$

$\langle \widetilde{\text{A-Cont}}, \circ, \leq, \widetilde{\tau}, \widetilde{\perp} \rangle$  can be shown to be a complete lattice. The function :

$$\text{Int} : \text{Arcs}^0 \times \widetilde{\text{A-Cont}} \rightarrow \text{A-Cont}$$

defines the interpretation of basic instructions. If  $\{\text{C}(q) \mid q \in \underline{\text{a-pred}}(n)\}$  is the set of input contexts of node  $n$ , then the output context on exit arc  $r$  of  $n$  ( $r \in \underline{\text{a-succ}}(n)$ ) is equal to  $\text{Int}(r, \text{C})$ .  $\text{Int}$  is supposed to be order-preserving :

$$\forall a \in \text{Arcs}, \forall (\text{Cv}', \text{Cv}'') \in \widetilde{\text{A-Cont}}^2,$$

$$\{\text{Cv}' \leq \text{Cv}''\} \Rightarrow \{\text{Int}(a, \text{Cv}') \leq \text{Int}(a, \text{Cv}'')\}$$

The local interpretation of elementary program constructs which is defined by  $\text{Int}$  is used to associate a system of equations with the program. We define

$$\widetilde{\text{Int}} : \widetilde{\text{A-Cont}} \rightarrow \widetilde{\text{A-Cont}} \mid \widetilde{\text{Int}}(\text{Cv}) = \lambda r. \text{Int}(r, \text{Cv})$$

It is easy to show that  $\widetilde{\text{Int}}$  is order-preserving. Hence it has fixpoints, [Tarski[55]]. Therefore the context vector resulting from the abstract interpretation  $\text{I}$  of program  $P$ , which defines the global properties of  $P$ , may be chosen to be one of the extreme solutions to the system of equations  $\text{Cv} = \widetilde{\text{Int}}(\text{Cv})$ .

## 5.2 Typology of Abstract Interpretations

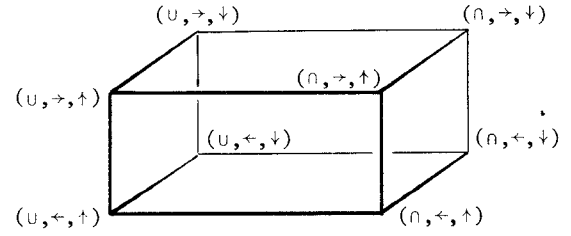
The restriction that "A-Cont" must be a complete semi-lattice is not drastic since Mac Neille[37] showed that any partly ordered set  $S$  can be embedded in a complete lattice so that inclusion is preserved, together with all greatest lower bounds and lowest upper bounds existing in  $S$ . Hence in practice the set of abstract contexts will be a lattice, which can be considered as a join ( $\cup$ ) semi-lattice or a meet ( $\cap$ ) semi-lattice, thus giving rise to two dual abstract interpretations.

It is a pure coincidence that in most examples (see 5.3.2) the  $\cap$  or  $\cup$  operator represents the effect of path converging. The real need for this operator is to define completeness which ensures  $\widetilde{\text{Int}}$  to have extreme fixpoints (see 8.4).

The result of an abstract interpretation was defined as a solution to forward ( $\rightarrow$ ) equations : the output contexts on exit arcs of node  $n$  are defined as a function of the input contexts on entry arcs of node  $n$ . One can as well consider a system of backward ( $\leftarrow$ ) equations : a context may be related to its successors. Both systems ( $\leftarrow$ ,  $\rightarrow$ ) may also be combined.

Finally we usually consider a maximal ( $\uparrow$ ) or minimal ( $\downarrow$ ) solution to the system of equations, (by agreement, maximal and minimal are related to the ordering  $\leq$  defined by  $(x \leq y) \iff (x \cup y = y) \iff (x \cap y = x)$ ). However known examples such as Manna and Shamir[75] show that the suitable solution may be somewhere between the extreme ones.

These choices give rise to the following types of abstract interpretations :



Examples :

Kildall[73] uses  $(n, \rightarrow, \uparrow)$ , Wegbreit[75] uses  $(u, \rightarrow, \uparrow)$ . Tenenbaum[74] uses both  $(u, \rightarrow, \uparrow)$  and  $(n, \leftarrow, \uparrow)$ .

## 5.3 Examples

### 5.3.1 Static Semantics of Programs

The static semantics of programs we defined in section 4 is an abstract interpretation :

$$\text{I}_{\text{SS}} = \langle \text{Contexts}, \cup, \subseteq, \text{Env}, \emptyset, \underline{\text{n-context}} \rangle$$

where Contexts,  $\cup$ ,  $\subseteq$ , Env,  $\emptyset$ ,  $\underline{\text{n-context}}$ , Context-Vectors,  $\widetilde{\cup}$ ,  $\widetilde{\subseteq}$ ,  $\widetilde{\text{F-Cont}}$  respectively correspond to A-Cont,  $\circ$ ,  $\leq$ ,  $\tau$ ,  $\perp$ ,  $\text{Int}$ , A-Cont,  $\circ$ ,  $\leq$ ,  $\widetilde{\text{Int}}$ .

### 5.3.2 Data Flow Analysis

Data flow analysis problems (see references in Ullman[75]) may be formalized as abstract interpretations of programs.

"Available expressions" give a classical example. An expression is available on arc  $r$ , if whenever control reaches  $r$ , the value of the expression has been previously computed, and since the last computation of the expression, no argument of the expression has had its value changed.

Let  $\text{Expr}_P$  be the set of expressions occurring in a program  $P$ . Abstract contexts will be sets of available expressions, represented by boolean vectors :

$$\text{B-vect} : \text{Expr}_P \rightarrow \{\text{true}, \text{false}\}$$

B-vect is clearly a complete boolean lattice. The interpretation of basic nodes is defined by :

$$\begin{aligned} \text{avail}(r, \text{Bv}) \\ \text{let } n \text{ be origin}(r) \text{ within} \\ \text{case } n \text{ in} \\ \quad \text{Entries} \Rightarrow \lambda e. \text{false} \\ \quad \text{Assignments } \cup \text{ Tests } \cup \text{ Junctions} \Rightarrow \\ \quad \lambda e. (\underline{\text{generated}}(n)(e) \text{ or } ((\text{and } \text{Bv}(p)(e)) \\ \quad \quad \text{p} \in \underline{\text{a-pred}}(n)) \\ \quad \quad \text{and } \underline{\text{transparent}}(n)(e))) \\ \text{esac} \end{aligned}$$

(Nothing is available on entry arcs. An expression  $e$  is available on arc  $r$  (exit of node  $n$ ) if either the expression  $e$  is generated by  $n$  or for all predecessors  $p$  of  $n$ ,  $e$  is available on  $p$  and  $n$  does not modify arguments of  $e$ ).

The available expressions are determined by the maximal solution (for ordering  $\lambda e. \text{false} \leq \lambda e. \text{true}$ ) of the system of equations :

$$\text{Bv} = \widetilde{\text{avail}}(\text{Bv})$$

The determination of available expressions, back-dominators, intervals, ... requires a forward system of equations. Some global flow problems, notably the live variables and very busy expressions require propagating information backward through the program graph, they are examples of backward systems of equations.

### 5.3.3 Remarks

Our formal definition of abstract interpretations has the completeness property since the model ensures the existence of a particular solution to the system of equations and therefore defines at least some global property of the program. It must also have the consistency property, that is define only correct properties of programs.

One can distinguish between syntactic and semantic abstract interpretations of a program. Syntactic interpretations are proved to be correct by reference to the program syntax (e.g. the algorithm for finding available expressions is justified by reasoning on paths of the program graph). By contrast semantic abstract interpretations must be proved to be consistent with the formal semantics of the language (e.g. constant propagation).

## 6. Consistent Abstract Interpretations

An "abstract" interpretation  $\bar{I} = \langle \bar{A}\text{-Cont}, \bar{\alpha}, \bar{\gamma}, \bar{\tau}, \bar{\iota}, \bar{\text{Int}} \rangle$  of a program is consistent with a "concrete" interpretation  $I = \langle C\text{-Cont}, \alpha, \gamma, \tau, \iota, \text{Int} \rangle$  if the context vector  $\bar{Cv}$  resulting from  $\bar{I}$  is a correct approximation of the context vector  $Cv$  resulting from the more refined interpretation  $I$ . This may be rigorously defined by establishing a correspondence ( $\alpha$  : abstraction) between concrete and abstract context vectors, and inversely ( $\gamma$  : concretization), and requiring :

$$6.0 \quad \{Cv \lesssim \tilde{\gamma}(\bar{Cv})\} \text{ and } \{\tilde{\alpha}(Cv) \lesssim \bar{Cv}\}$$

In words the abstract context vector must at least contain the concrete one, (but not only the concrete one).

If  $f : D \rightarrow D'$  we note  $\tilde{D} = \text{Arcs}^0 \rightarrow D$  and  $\tilde{D}' = \text{Arcs}^0 \rightarrow D'$  and  $\tilde{f} : \tilde{D} \rightarrow \tilde{D}' = \lambda d. (\lambda r. f(d(r)))$ .

We will suppose  $\alpha$  and  $\gamma$  to satisfy the following hypothesis :

$$6.1 \quad \alpha : C\text{-Cont} \rightarrow \bar{A}\text{-Cont}, \quad \gamma : \bar{A}\text{-Cont} \rightarrow C\text{-Cont}$$

$$6.2 \quad \alpha \text{ and } \gamma \text{ are order-preserving}$$

$$6.3 \quad \forall \bar{x} \in \bar{A}\text{-Cont}, \quad \bar{x} = \alpha(\gamma(\bar{x}))$$

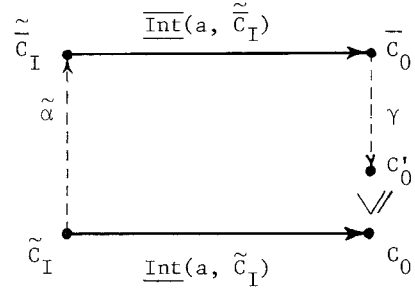
$$6.4 \quad \forall x \in C\text{-Cont}, \quad x \leq \gamma(\alpha(x))$$

Intuitively, hypothesis 6.2 is necessary because context inclusion (that is property comparison) must be preserved by the abstraction or concretization process. 6.3 requires that concretization introduces no loss of information. It implies that  $\alpha$  is surjective and  $\gamma$  is injective. 6.4 introduces the idea of approximation : the abstraction  $\alpha(C)$  of a concrete context  $C$  may introduce some loss of information so that when concretizing again  $\gamma(\alpha(C))$  we may get a larger context  $\gamma(\alpha(C)) \geq C$ . Note that it is easy to prove properties corresponding to 6.1-6.4 for  $\alpha$  and  $\gamma$ .

Instead of the global hypothesis 6.0 we will use the following local hypothesis on the concrete and abstract interpretations of primitive language constructs :

$$6.5 \quad \text{and} \quad \begin{aligned} & \{ \forall (a, \bar{x}) \in \text{Arcs} \times \bar{A}\text{-Cont}, \\ & \quad \gamma(\bar{\text{Int}}(a, \bar{x})) \geq \text{Int}(a, \gamma(\bar{x})) \} \\ & \{ \forall (a, x) \in \text{Arcs} \times C\text{-Cont}, \\ & \quad \bar{\text{Int}}(a, \alpha(x)) \geq \alpha(\text{Int}(a, x)) \} \end{aligned}$$

These two hypothesis are in fact equivalent (lemma L2 in appendix 12). The following schema illustrates 6.5, i.e. the idea of abstract simulation of concrete computations :



Suppose we want to compute the concrete output context  $C_0$  (associated with arc  $a$ ) resulting from concrete input contexts  $\tilde{C}_I : C_0 = \text{Int}(a, \tilde{C}_I)$ . We can as well approximate this computation in the abstract universe, and get  $\bar{C}'_0 = \gamma(\bar{\text{Int}}(a, \tilde{\alpha}(\tilde{C}_I)))$ . 6.5 requires  $\bar{C}'_0$  to contain at least  $C_0$ , that is  $C_0 \leq \bar{C}'_0$ . On the contrary we do not require  $\bar{C}'_0$  to contain at most  $C_0$ , that is  $\bar{C}'_0 \leq C_0$  is not compulsory.

We will say that  $I$  is a refinement of  $\bar{I}$ , or that  $\bar{I}$  is an abstraction of  $I$ , denoted  $I \leq (\alpha, \gamma)\bar{I}$ , if and only if there exist  $\alpha$  and  $\gamma$  satisfying hypothesis 6.1 to 6.5.

Note that  $I \leq (\alpha, \gamma)\bar{I}$  imposes a local consistency of the interpretations  $I$  and  $\bar{I}$ , at the level of primitive language constructs (6.5). Theorems T1 and T2 of Appendix 12 then prove 6.0 which defines the global consistency of  $I$  and  $\bar{I}$  at the program level.

In particular if we take

$$I_{SS} = \langle \text{Contexts}, \cup, \subseteq, \text{Env}, \emptyset, \text{n-context} \rangle$$

any abstract interpretation  $\bar{I}$  of  $P$ , consistent with  $I_{SS}$  ( $I_{SS} \leq (\alpha, \gamma)\bar{I}$ ) is consistent with the semantics of  $P$ , which implies :

$\forall q \in \text{Arcs}$ , let  $\bar{Cv}$  be the result of  $\bar{I}$ ,

$$\{ \exists n \in \mathbb{N}, \exists i_s \in \text{I-states} \mid \langle q, e \rangle = \text{n-state}^n(i_s) \} \\ \implies \{ e \in \gamma(\bar{Cv}(q)) \}$$

As previously noticed, the abstract interpretations will not in general be powerful enough to establish the reciprocal.

### Example : Deductive Semantics of Programs

Contexts will be predicates such as  $P(x_1, \dots, x_n) \in \text{Pred}$  over the program variables  $(x_1, \dots, x_n) \in \text{Ident}^n$  which are the free variables in the predicate. The abstract interpretation is then :

$$I_{DS} = \langle \text{Pred}, \text{or}, \implies, \text{true}, \text{false}, \text{n-pred} \rangle$$

where  $\underline{n\text{-pred}}$  defines Floyd[67]'s strongest post condition :

$$\begin{aligned} \underline{n\text{-pred}}(r, \underline{Pv}) = & \underline{\text{let}}(n \text{ be } \underline{\text{origin}}(r)), (p \text{ be } \underline{a\text{-pred}}(\underline{\text{origin}}(r))) \underline{\text{within}} \\ & \underline{\text{case } n \text{ in}} \\ & \quad \underline{\text{Entries}} \implies (\forall x \in \underline{\text{Ident}}, x = \perp_{\underline{\text{Values}}}) \\ & \quad \underline{\text{Junctions}} \implies \text{or} \quad (\underline{Pv}(q)) \\ & \quad \quad q \in \underline{a\text{-pred}}(n) \\ & \quad \underline{\text{Tests}} \implies \underline{\text{case } r \text{ in}} \\ & \quad \quad \{ \underline{a\text{-succ-t}}(n) \} \implies \underline{Pv}(p) \text{ and } \underline{\text{test}}(n) \\ & \quad \quad \{ \underline{a\text{-succ-f}}(n) \} \implies \underline{Pv}(p) \text{ and } \underline{\text{not test}}(n) \\ & \quad \underline{\text{esac}} \\ & \quad \underline{\text{Assignments}} \implies \\ & \quad \quad \underline{\text{let}} \ (P \text{ be } \underline{Pv}(p)), (x \text{ be } \underline{\text{id}}(n)), \\ & \quad \quad \quad (\underline{e} \text{ be } \underline{\text{expr}}(n)) \underline{\text{within}} \\ & \quad \quad \quad (\exists v \in \underline{\text{Values}} \mid P[v/x] \text{ and } x = e[v/x]) \\ & \quad \underline{\text{esac}} \end{aligned}$$

The "invariants" of the program are defined by the least fixpoint of  $\underline{n\text{-pred}}$  (least for ordering  $\lesssim$  ( $\approx$ )), so that an invariant implies any other correct assertion).

The deductive semantics is easily validated by proving that  $I_{SS} \leq (\alpha, \gamma) I_{DS}$ , where :

$$\begin{aligned} \alpha : \text{Contexts} &\rightarrow \text{Pred} \\ &= \lambda C. (\text{or } ( \text{and } (x = e(x)) \\ &\quad e \in C \quad x \in \underline{\text{Ident}} \\ \gamma : \text{Pred} &\rightarrow \text{Contexts} \\ &= \lambda P. \{ e \mid P[e(x)/x, x \in \underline{\text{Ident}}] \} \end{aligned}$$

The main point is to justify Hoare[67]'s proof rules by showing :

$$\{ \forall a \in \text{Arcs}, \forall \underline{Pv} \in \underline{\text{Pred}}, \\ \alpha(\underline{n\text{-context}}(a, \gamma(\underline{Pv}))) \implies \underline{n\text{-pred}}(a, \underline{Pv}) \}$$

See Hoare and Lauer[74], Ligler[75]. In particular Ligler[75] shows clearly that the proof can be done only when considering realizable Contexts and programs involving "clean" basic constructs (e.g. constructs excluding non-termination, errors, side-effects, sharing between identifiers, ...).

Once  $I_{SS} \leq (\alpha, \beta) I_{DS}$  has been proved, we know that the deductive semantics gives a valid proof technique, which will never permit a false theorem to be deduced :

$$\begin{aligned} \forall q \in \text{Arcs}, \text{ let } \underline{Pv} \text{ be the result of } I_{DS}, \\ \{ \exists n \geq 0, \exists i_s \in \text{I-states} \mid \langle q, e \rangle = \underline{n\text{-state}}^n(i_s) \} \\ \implies \{ \underline{Pv}(q) \implies \alpha(e) \} \end{aligned}$$

## 7. The Lattice of Abstract Interpretations

The relation  $\leq$  comparing the levels of abstraction of two interpretations is a quasi-ordering since it is :

$$\begin{aligned} \text{reflexive : } (I \leq (\gamma, \gamma)I) \text{ where } \gamma = \lambda x. x \text{ is} \\ \text{the identify function,} \\ \text{transitive : } (I \leq (\alpha_1, \gamma_1)I') \text{ and} \\ (I' \leq (\alpha_2, \gamma_2)I'') \text{ imply} \\ I \leq (\alpha_1 \circ \alpha_2, \gamma_2 \circ \gamma_1)I''. \end{aligned}$$

The relation  $\equiv$  on abstract interpretations defined by :

$$\{ I \equiv I' \} \iff \{ (I \leq I') \text{ and } (I' \leq I) \}$$

is an equivalence relation. We have :

$$\{ I \equiv (\beta)I' \} \iff \{ \beta \text{ is an isomorphism between the algebras } I \text{ and } I' \}$$

The proof gives some insight in the abstraction process :

$$1 - \{ I \equiv (\beta)I' \} \implies \{ (I \leq (\beta, \beta^{-1})I') \text{ and } (I' \leq (\beta^{-1}, \beta)I) \}$$

2 - reciprocally,

If  $I \leq (\alpha_1, \gamma_1)I'$ , let  $\equiv (\alpha_1)$  be the equivalence relation defined on  $I$  (properly speaking, on the set of abstract contexts of  $I$ ) by :

$$\{ x \equiv (\alpha_1)y \} \iff \{ \alpha_1(x) = \alpha_1(y) \}$$

$\forall x' \in I'$ , each equivalence class  $C_x = \{ x \in I \mid \alpha_1(x) = x' \}$  has a least upper bound  $x$  which is  $\gamma_1(x')$ . Hence the projection  $\alpha_1 \mid \gamma_1(I')$  of  $\alpha_1$  on  $\gamma_1(I')$  is a bijection from the set  $\gamma_1(I')$  of representers of the equivalence classes on  $I$ . Let us show now that under the hypothesis  $I \leq (\alpha_1, \gamma_1)I'$  and  $I' \leq (\alpha_2, \gamma_2)I$ ,  $\alpha_1$  is bijective.

$\alpha_1 \mid \gamma_1(I')$  and  $\alpha_2 \mid \gamma_2(I)$  are bijections, hence  $\forall x' \in I'$ ,  $\exists! x$  (unique)  $\in \gamma_1(I')$  such that  $x' = (\alpha_1 \mid \gamma_1(I'))(x)$ . Likewise,  $x \in \gamma_1(I') \implies x \in I \implies \exists! x'' \in \gamma_2(I) \mid x = (\alpha_2 \mid \gamma_2(I))(x'')$ .

Therefore,  $\forall x' \in I'$ ,  $\exists! x'' \in \gamma_2(I) \mid x' = (\alpha_1 \mid \gamma_1(I')) \circ (\alpha_2 \mid \gamma_2(I))(x'')$ . Thus  $(\alpha_1 \mid \gamma_1(I')) \circ (\alpha_2 \mid \gamma_2(I))$  is a bijection between  $\gamma_2(I)$  and  $I'$ . Since  $(\alpha_2 \mid \gamma_2(I))^{-1}$  is a bijection between  $I$  and  $\gamma_2(I)$ , the composition

$$\begin{aligned} (\alpha_1 \mid \gamma_1(I')) \circ (\alpha_2 \mid \gamma_2(I)) \circ (\alpha_2 \mid \gamma_2(I))^{-1} \\ = (\alpha_1 \mid \gamma_1(I')) \end{aligned}$$

is a bijection between  $I$  and  $I'$ , hence  $\alpha_1$  is a bijection between  $I$  and  $I'$  which is trivially an algebraic morphism. ( $\alpha_1$  is isotone, its inverse  $\alpha_1^{-1} = \gamma_1$  is isotone and  $\alpha_1(\underline{\text{Int}}(a, X)) = \underline{\text{Int}}'(a, \tilde{\alpha}_1(X))$  Q.E.D.

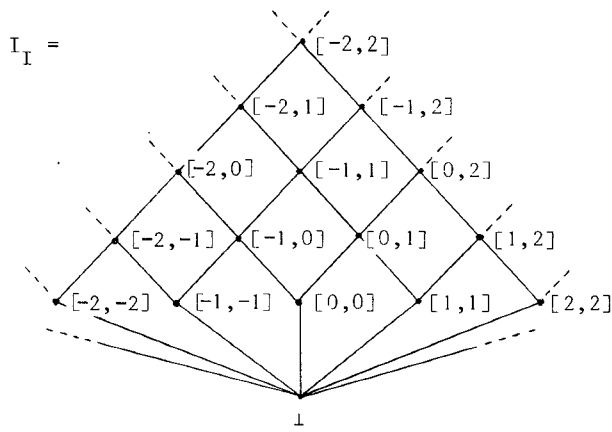
Let  $I$  be the set of abstract interpretations of a program, if equivalent interpretations are identified, the quasi-ordering  $\leq$  becomes a partial ordering.

In particular, we can restrict  $I$  to be set of interpretations which abstract  $I_{SS}$ .  $I$  is then a lattice, (with ordering  $\leq$ ) which is isomorphic with a subset of the lattice of equivalence relations on Contexts.

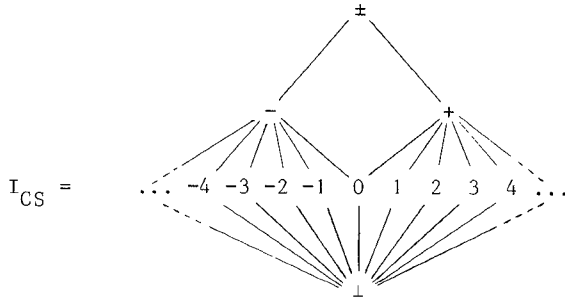
Example :

Let  $P$  be a program with a single integer variable, (the generalization is obvious). Environments will be integers (the value of the variable). Contexts are sets of integers (the set of values at some program point).

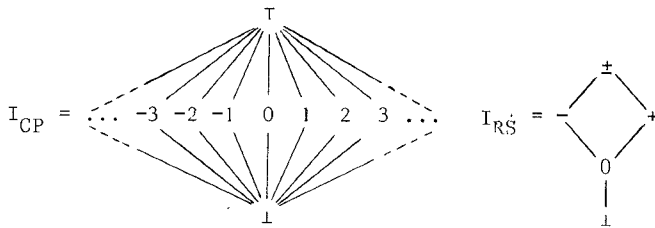
A context  $S$  may be abstracted by a closed interval  $\alpha(S) = [\min(S), \max(S)]$ . When  $S$  is infinite the bounds will eventually be  $-\infty$  or  $+\infty$ .  $\gamma([a, b]) = \{ x \mid a \leq x \leq b \}$ . The abstract contexts are then, (Cousot[76]) :



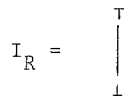
A further abstraction may be :  
 $\alpha([a, b]) = \text{if } a = b \text{ then } a \text{ elseif } a \geq 0 \text{ then } +$   
 $\text{elseif } b \leq 0 \text{ then } - \text{ else } \pm \text{ fi. } \gamma(n) = [n, n],$   
 $\gamma(+) = [0, +\infty], \gamma(-) = [-\infty, 0], \gamma(\pm) = [-\infty, +\infty].$   
 The abstract contexts are then :



This interpretation may be abstracted by two non-comparable abstractions :



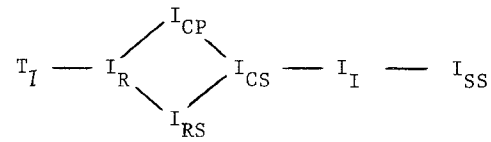
$I_{CP}$  is used by Kildall[73] for constant propagation.  $I_{RS}$  might be used to apply the rules of signs. Both interpretations may be abstracted by :



which may be used to check that any vertex in the program graph is reachable from the entry nodes. Finally, the most abstract interpretation is the upper bound of  $I$  :

$$T_I = \langle \{I\}, \lambda(x, y). I, t, I, I, \lambda(a, C). I \rangle$$

where  $t$  is the relation which is always true. We have exhibited a sublattice of  $I$  which is :



## 8. Abstract Evaluation of Programs

The system of equations :

$$Cv : \widetilde{\text{Int}}(Cv)$$

resulting from an interpretation  $I = \langle A\text{-Cont}, \circ, \leq, \tau, \perp, \text{Int} \rangle$  of a program  $P$  may be solved by "elimination" methods, (e.g. Tarjan[75]). Otherwise, one can use an "iterative" algorithm which computes Kleene's sequence (L4 of Appendix 12) :

$$Cv := (C := \widetilde{1}; \text{until } C = \widetilde{\text{Int}}(C) \text{ do } C := \widetilde{\text{Int}}(C) \text{ repeat}; C)$$

### 8.1 Correctness

If  $\text{Int}$  is supposed to be a complete morphism (i.e. infinitely distributive over  $\circ$ ) then  $Cv$  is the least fixpoint of  $\widetilde{\text{Int}}$ . (e.g. Kildall[73], since in a semi-lattice of finite length, any distributive function is a complete morphism). Under the weaker assumption that  $\text{Int}$  is continuous, the limit  $Cv$  of Kleene's sequence can also be shown to be the least fixpoint of  $\widetilde{\text{Int}}$  (e.g. Wegbreit[75], since in a well-founded semi-lattice, any isotone function is continuous). Finally, if  $\text{Int}$  is only supposed to be isotone,  $Cv$  is an approximation ( $\leq$ ) of the least fixpoint (e.g. Kam and Ullman[75]).

### 8.2 Termination

The abstract evaluation terminates iff Kleene's sequence is finite. This may be the case because  $A\text{-Cont}$  is finite (e.g. type checking in ALGOL 60, Naur[65]), or a finite subset only is to be considered for any particular program (e.g. type checking in ALGOL 68), or  $A\text{-Cont}$  may be of finite length  $m$  (the length of any strictly increasing chain is bounded by  $m$ , Kildall[73], Wegbreit[75]) or  $A\text{-Cont}$  may satisfy the ascending chain condition (every strictly increasing chain is finite, although not bounded). A lattice may have infinite chains, although  $\text{Int}$  is chosen so that Kleene's sequences are finite. Finally an infinite Kleene's sequence may be arbitrarily truncated (to get a lower bound of its limit), some induction principle (Sintzoff[75]) or heuristics (Katz and Manna[76]) may be used to pass to the limit, or approximate it, (Cousot[76]).

### 8.3 Efficiency

In practice efficient versions of the Kleene's sequence are used. These consist in a symbolic execution of the program which propagates information along paths of the program until stabilization. A specification of order of information propagation may lead to optimal algorithms for specific applications (references in Tarjan[76]).

#### 8.4 Example : Performance Analysis of Programs

The performance of programs may be analyzed by deriving for each program point the final value of an imaginary counter which is incremented each time control goes through that point.

Let  $A\text{-Cont}$  be the lattice  $\mathbb{R}^+$  of positive real numbers augmented by the upper bound  $\infty$ , with natural ordering  $\leq$ . The abstract interpretation :

$$I_P = \langle \mathbb{R}^+, \max, \leq, 0, \infty, \text{Kir} \rangle$$

may be used to derive the mean values of the counters using Kirchhoff's law of conservation of flow :

```
Kir(r, Cv) =
  let n be origin(r) within
  case n in
    Entries ==> 1 {unique entry node}
    Junctions u Assignments ==> \sum_{p \in a\text{-pred}(n)} Cv(p)
    Tests ==>
      case r in
        {a-succ-t(n)} ==> Cv(a-pred(n)) * Prob(test(n) = true)
        {a-succ-f(n)} ==> Cv(a-pred(n)) * (1 - Prob(test(n) = true))
      esac
    esac
```

The main difficulty is to obtain the probability  $\text{Prob}(\text{test}(n) = \text{true})$  of taking the true path at a test node  $n$ . Suppose the values of these probabilities can be determined (from hypothesis on the input data).

For fixed probabilities, the function  $\widetilde{\text{Kir}}$  is clearly continuous (although it is not a complete morphism) since

$$\text{if } \underline{Cv}_0 \lesssim \underline{Cv}_1 \lesssim \dots \lesssim \underline{Cv}_n \lesssim \dots$$

$$\text{then } \max_{i=0}^{\infty} \left( \sum_{p \in a\text{-pred}(n)} \underline{Cv}_i(p) \right) = \sum_{p \in a\text{-pred}(n)} \left( \max_{i=0}^{\infty} \underline{Cv}_i(p) \right)$$

$$\text{and } \max_{i \in \Delta} (n_i * q) = \left( \max_{i \in \Delta} n_i \right) * q.$$

The least fixpoint of  $\widetilde{\text{Kir}}$  is the limit of Kleene's sequence (the length of the sequence is in general infinite) :

- Let  $P$  be the program "begin  $L$  : go to  $L$  end". The number  $n$  of iterations in the loop is given by the minimal solution to the equation  $n = n + 1$  which is  $\infty$  limit of  $0 + 1 + 1 + 1 + \dots$

- Let  $P$  be the program "begin while  $T$  do  $I$  end". The number  $n$  of times the expression  $T$  is tested is given by the minimal solution to the equation  $n = 1 + q * n$  where  $q$  is the probability of  $T$  to be true.  $n$  may be determined by the limit of Kleene's sequence :

$$0 + 1 + q + q^2 + \dots + q^l + \dots$$

which is an infinite series. Its sum is  $\frac{1}{1-q}$ .

This abstract interpretation leads to a system of linear equations. Kleene's sequence corresponds to the Jacobi's iterative method (for numerical coefficients).

#### 9. Fixpoints Approximation Methods

When the extreme fixpoints of the system of equations established for an abstract interpretation  $I$  of a program  $P$  cannot be computed in finitely many steps, they can be approximated. A more abstract interpretation  $\bar{I}$  ( $I \leq \bar{I}$ ) may be used for that purpose (e.g. Tenenbaum[74]). It is often better to make approximations in  $I$ , for example by "accelerating the convergence" of Kleene's sequences.

##### 9.1 Finite Iterative and Increasing Approximation of the Least Fixpoint Starting from a Lower Bound

Let  $I = \langle A\text{-Cont}, \circ, \leq, 1, \tau, \text{Int} \rangle$  be an interpretation of  $P$ . When the least fixpoint  $\underline{Cv}$  of  $\text{Int}$  is unreachable, we look for an upper bound  $\underline{Ub}$  of  $\underline{Cv}$ , since according to the correctness requirement 6.0,  $\underline{Cv} \lesssim \widetilde{\gamma}(\underline{Cv})$  and  $\underline{Cv} \lesssim \underline{Ub}$  implies  $\underline{Cv} \lesssim \widetilde{\gamma}(\underline{Ub})$ .

##### 9.1.1 Increasing Approximation Sequence

Let  $\widetilde{A\text{-int}} : A\text{-Cont} \rightarrow A\text{-Cont}$  be such that :

$$9.1.1.1 \quad \{ \forall n \geq 0, C = \widetilde{A\text{-int}}^n(\tilde{I}) \text{ and } \text{not}(\widetilde{\text{Int}}(C) \lesssim C) \} \\ \Rightarrow \{ C \lesssim \widetilde{\text{Int}}(C) \lesssim \widetilde{A\text{-int}}(C) \}.$$

$$9.1.1.2 \quad \text{Every infinite sequence } \tilde{I}, \widetilde{A\text{-int}}(\tilde{I}), \dots, \widetilde{A\text{-int}}^n(\tilde{I}), \dots \text{ is not strictly increasing.}$$

The approximation sequence  $S_0, \dots, S_n, \dots$  is recursively defined by :

$$9.1.1.3 \quad S_0 = \tilde{I} \\ S_{n+1} = \text{if } \text{not}(\widetilde{\text{Int}}(S_n) \lesssim S_n) \text{ then } \widetilde{A\text{-int}}(S_n) \\ \text{else } S_n \\ \text{fi}$$

We now prove that  $\exists m$  finite such that :

$$S_0 \lesssim S_1 \lesssim \dots \lesssim S_m = S_{m+1} = \dots$$

Let  $m$  be the least natural number (eventually infinite) such that  $S_m = S_{m+1}$ .  $\forall k \in [0, m[$ , we know from 9.1.1.3 that  $\text{not}(\widetilde{\text{Int}}(S_k) \lesssim S_k)$ . Whence by definition of the ordering  $\lesssim$ ,  $S_k \neq \widetilde{\text{Int}}(S_k) \gtrsim S_k$ .

Since  $S_k \lesssim \widetilde{\text{Int}}(S_k) \gtrsim S_k$  is always true, we can state that  $S_k \gtrsim \widetilde{\text{Int}}(S_k) \gtrsim S_k$ . Besides  $\text{not}(\widetilde{\text{Int}}(S_k) \lesssim S_k)$  and 9.1.1.1 imply :

$$S_{k+1} = \widetilde{A\text{-int}}(S_k) \gtrsim \widetilde{\text{Int}}(S_k) \gtrsim S_k$$

and therefore we conclude  $S_{k+1} \gtrsim S_k$ ,  $\forall k \in [1, m[$ . Moreover 9.1.1.2 implies that  $m$  is finite. Q.E.D.

Let  $\underline{Cv}$  be the least fixpoint of  $\widetilde{\text{Int}}$ , it is the greatest lower bound of the set of  $X \in A\text{-Cont}$  such that  $\widetilde{\text{Int}}(X) \gtrsim X$  (Tarski[55]) hence :

$$\forall X \in A\text{-Cont}, \{ \widetilde{\text{Int}}(X) \gtrsim X \} \Rightarrow \{ \underline{Cv} \gtrsim X \}$$

Since  $S_m = S_{m+1}$  we have  $\widetilde{\text{Int}}(S_m) \gtrsim S_m$  and therefore  $\underline{Cv} \gtrsim S_m$ .  $S_m$  is a correct approximation of  $\underline{Cv}$ .



### 9.1.2 Generalization of Kleene's Ascending Sequence

When A-Cont satisfies the ascending chain condition one can choose  $\widetilde{\text{A-int}}$  to be  $\text{Int}$  and therefore the approximation sequence generalizes Kleene's sequence and the related methods.

### 9.1.3 Widening in Increasing Approximation sequences

The definition of the approximate interpretation  $\widetilde{\text{A-int}}$  in 9.1.1 is global. We now indicate a way to construct  $\widetilde{\text{A-int}}$  by local modifications to  $\text{Int}$ .

Let  $(q, r) \in \text{Arcs}^2$ , we say that the context associated to  $q$  is dependent on the context associated to  $r$ , if and only if :

$$\{\exists \text{Cv} \in \widetilde{\text{A-Cont}}, \exists \text{C} \in \text{A-Cont} \mid \text{Int}(q, \text{Cv}) \neq \text{Int}(q, \text{C}[r])\}$$

(e.g. in a forward system of equations the context associated to  $q$  may only depend on the contexts associated with the immediate predecessor arcs of  $q$ ). In the system of equations  $\text{Cv} = \widetilde{\text{Int}}(\text{Cv})$  we define a cycle to be a sequence  $\langle q_1, \dots, q_n \rangle$  of arcs, such that  $\forall i \in [1, n], \text{Cv}(q_{i+1})$  depends on  $\text{Cv}(q_i)$  and  $\text{Cv}(q_1)$  depends on  $\text{Cv}(q_n)$ . (e.g. in a forward interpretation a cycle corresponds to a loop in the program).

In any infinite strictly increasing Kleene's sequence  $\text{Cv}_1, \dots, \text{Cv}_m, \dots$  since Arcs is finite there is some arc  $q$  for which the sequence  $\text{Cv}_1(q), \dots, \text{Cv}_m(q), \dots$  never stabilizes. Therefore  $q$  must belong to a cycle or the contexts associated to  $q$  transitively depend on the contexts associated to some other arc  $r$  which itself belongs to a cycle. The sequence of contexts associated to any arc of that cycle never stabilizes. In order to avoid this phenomenon, we introduce :

- The binary operation  $\nabla$  called widening defined by :
  - 9.1.3.1  $\nabla : \text{A-Cont} \times \text{A-Cont} \rightarrow \text{A-Cont}$
  - 9.1.3.2  $\forall (C, C') \in \text{A-Cont}^2, C \circ C' \leq C \nabla C'$
  - 9.1.3.3 Every infinite sequence  $s_0, \dots, s_n, \dots$  of the form  $s_0 = C_0, \dots, s_n = s_{n-1} \nabla C_n, \dots$  (where  $C_0, \dots, C_n, \dots$  are arbitrary abstract contexts) is not strictly increasing.
- The set W-arcs of widening arcs, which is one of the minimal sets of arcs such that any cycle  $\langle q_1, \dots, q_n \rangle$  of the system of equations  $\text{Cv} = \widetilde{\text{Int}}(\text{Cv})$  contains at least a widening arc :  $\exists i \in [1, n] \mid q_i \in \text{W-arcs}$ . (e.g. in a forward interpretation on a reducible program graph, W-arcs may be chosen to be the set of exit arcs of the junction nodes which are interval headers. On irreducible graphs an arbitrary choice has to be made so that any loop of the program goes through a widening arc).
- The approximate interpretation  $\widetilde{\text{A-int}} : \text{Arcs}^0 \times \text{A-Cont} \rightarrow \text{A-Cont}$  defined by :
  - 9.1.3.4  $\widetilde{\text{A-int}} = \lambda(q, \text{Cv}). \text{ if } q \in \text{W-arcs} \text{ then } \frac{\text{Cv}(q) \nabla \text{Int}(q, \text{Cv})}{\text{else}} \frac{\text{Int}(q, \text{Cv})}{\text{fi}}$

As before, we define :

$$9.1.3.5 \quad \widetilde{\text{A-int}} = \lambda \text{Cv}. (\lambda q. \widetilde{\text{A-int}}(q, \text{Cv}))$$

Now we have to show that this definition of  $\widetilde{\text{A-int}}$  satisfies the requirements 9.1.1.2 and 9.1.1.1.

Let us consider a sequence  $S_0 = \tilde{1}, \dots, S_{n+1} = \widetilde{\text{A-int}}(S_n), \dots$ . We show that this sequence is increasing that is to say :

$$9.1.3.6 \quad S_n \leq \widetilde{\text{A-int}}(S_n), \forall n \geq 0.$$

Trivially for  $n=0, S_0 = \tilde{1} \leq \widetilde{\text{A-int}}(S_0)$ . For the induction step, suppose the result to be true for  $n \leq m$ . Let us prove that :

$$\begin{aligned} S_{m+1} &\leq \widetilde{\text{A-int}}(S_{m+1}) \\ \iff S_{m+1}(q) &\leq \widetilde{\text{A-int}}(q, S_{m+1}), \forall q \in \text{Arcs}. \\ \text{If } q \in \text{W-arcs}, &\text{ then } \widetilde{\text{A-int}}(q, S_{m+1}) = S_{m+1}(q) \nabla \text{Int}(q, S_{m+1}) \geq S_{m+1}(q) \circ \text{Int}(q, S_{m+1}) \\ &\geq S_{m+1}(q). \\ \text{If } q \notin \text{W-arcs}, &\text{ then } \widetilde{\text{A-int}}(q, S_{m+1}) = \text{Int}(q, S_{m+1}) \\ &= \text{Int}(q, S_m) \leq \widetilde{\text{A-int}}(q, S_{m+1}) \text{ since } S_m \leq S_{m+1} \text{ and } \text{Int} \text{ is order preserving. Moreover from } q \notin \text{W-arcs} \text{ and } 9.1.3.4 \text{ we get } \text{Int}(q, S_m) = \widetilde{\text{A-int}}(q, S_m) \text{ and therefore } S_{m+1}(q) = \widetilde{\text{A-int}}(q, S_m) \leq \widetilde{\text{A-int}}(q, S_{m+1}). \end{aligned}$$

Finally  $S_{m+1} \leq \widetilde{\text{A-int}}(S_{m+1})$ , Q.E.D.

An infinite sequence  $S_0 = \tilde{1}, \dots, S_n = \widetilde{\text{A-int}}^n(\tilde{1}), \dots$  cannot be strictly increasing since otherwise there would exist some widening arc  $q$  for which the sequence  $S_0(q), \dots, S_n(q), \dots$  would never stabilize thus contradicting 9.1.3.3.

We now prove 9.1.1.1 that is to say that :

$$\begin{aligned} \forall n \geq 0, S_n &= \widetilde{\text{A-int}}^n(\tilde{1}) \\ \text{implies } S_n &\sim \widetilde{\text{Int}}(S_n) \leq \widetilde{\text{A-int}}(S_n) \end{aligned}$$

$$\begin{aligned} \iff (S_n \sim \widetilde{\text{Int}}(S_n))(q) &\leq \widetilde{\text{A-int}}(S_n)(q), \forall q \in \text{Arcs} \\ \iff S_n(q) \circ \text{Int}(q, S_n) &\leq \widetilde{\text{A-int}}(q, S_n) \text{ (see 9.1.3.5)} \end{aligned}$$

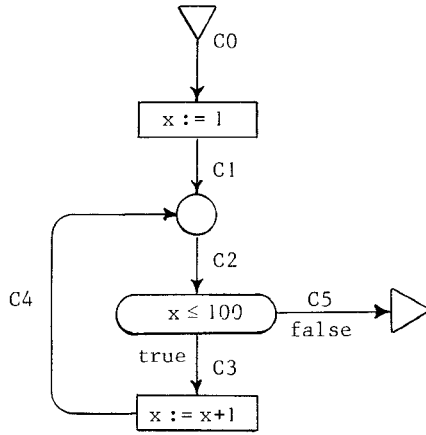
If  $q \in \text{W-arcs}$ , we have  $\widetilde{\text{A-int}}(q, S_n) = S_n(q) \nabla \text{Int}(q, S_n) \geq S_n(q) \circ \text{Int}(q, S_n)$  by 9.1.3.2. If now  $q \notin \text{W-arcs}$  we must show :

$$\begin{aligned} S_n(q) \circ \text{Int}(q, S_n) &\leq \text{Int}(q, S_n) \\ \iff S_n(q) \circ \text{Int}(q, S_n) &= \text{Int}(q, S_n) \\ \iff S_n(q) &\leq \text{Int}(q, S_n) \\ \iff S_n(q) &\leq \widetilde{\text{A-int}}(q, S_n) \text{ by 9.1.3.4} \end{aligned}$$

which is true, from 9.1.3.6, Q.E.D.

### 9.2 Example : Bounds of Integer Variables

In a PASCAL program operating on arrays, the compiler should ensure that arrays are subscripted only by indices within bounds. For that purpose one can use the lattice  $I_I$  of section 7. Let us take an obvious example :



Let us note  $[a, b]$  where  $a \leq b$  the predicate  $a \leq x \leq b$ . The system of equations corresponding to the example is :

- (0)  $C0 = [ , ]$
- (1)  $C1 = [1, 1]$
- (2)  $C2 = C1 \cup C4$
- (3)  $C3 = C2 \cap [-\infty, 100]$
- (4)  $C4 = C3 + [1, 1]$
- (5)  $C5 = C2 \cap [101, +\infty]$

Assignment statements are treated using an interval arithmetic (e.g.  $[i, j] + [k, l] = [i+k, j+l]$  naturally extended to include the case of the empty interval). Similarly tests are treated using an "interval logic". Since there exist infinite Kleene's sequences (e.g.  $[ , ] < [0, 0] < [0, 1] < \dots < [0, +\infty]$  for the program  $x := 0 ; \text{while true do } x := x+1$ ), we must use an approximation sequence. Hence the results will be somewhat inaccurate but runtime subscript tests may be inserted in the absence of certainty.

Let us define the widening  $\nabla$  of intervals by :

- $[ , ]$  is the null element of  $\nabla$
- $[i, j] \nabla [k, l] = \begin{cases} \text{if } k < i \text{ then } -\infty \text{ else } i \text{ fi,} \\ \text{if } l > j \text{ then } +\infty \text{ else } j \text{ fi} \end{cases}$

$\nabla$  satisfies the requirements of 9.1.3. According to 9.1.3.4 the system of equations is modified by :

$$(2) \quad C2 = C2 \nabla (C1 \cup C4)$$

The corresponding approximation sequence is :

- $C_i = [ , ]$  for  $i \in [0, 5]$
- \*  $C1 = [1, 1]$
- $C2 = C2 \nabla (C1 \cup C4)$
- $= [ , ] \nabla ([1, 1] \cup [ , ])$
- $= [ , ] \nabla [1, 1]$
- $= [1, 1]$
- $C3 = C2 \cap [-\infty, 100]$
- $= [1, 1] \cap [-\infty, 100]$
- $= [1, 1]$
- $C4 = C3 + [1, 1]$
- $= [1, 1] + [1, 1]$
- $= [2, 2]$
- $C2 = C2 \nabla (C1 \cup C4)$
- $= [1, 1] \nabla ([1, 1] \cup [2, 2])$
- $= [1, 1] \nabla [1, 2]$
- \*  $C2 = [1, +\infty]$
- $C3 = C2 \cap [-\infty, 100]$
- $= [1, +\infty] \cap [-\infty, 100]$

- \*  $C3 = [1, 100]$
- $C4 = C3 + [1, 1]$
- $= [1, 100] + [1, 1]$
- \*  $C4 = [2, 101]$
- Note :  $C1 \cup C4 = [1, 101] \leq C2 = [1, +\infty]$
- stop on that path.
- $C5 = C2 \cap [101, +\infty]$
- $= [1, +\infty] \cap [101, +\infty]$
- \*  $C5 = [101, +\infty]$
- exit, stop.

The final context on each arc is marked by a star \*. Note that the results are approximate ones, (e.g. C5).

In this example the widening is a very rough operation which introduces a great loss of information. However it can be seen in the trace that tests behave like filters. Furthermore, for PASCAL like languages, one can first use the bounds given in the declaration of  $x$  before widening to infinite limits.

### 9.3 Finite Iterative and Decreasing Approximation of the Least Fixpoint Starting from an Upper Bound

The ascending approximation sequence leads to an upper bound  $S_m = A\text{-int}^m(\gamma)$  of the least fixpoint  $\underline{Cv}$  of  $\underline{\text{Int}} : \underline{Cv} \lesssim S_m$ . Moreover  $\underline{\text{Int}}(S_m) \lesssim S_m$ . Since  $\underline{\text{Int}}$  is order preserving, this implies that :

$$S_m \gtrsim \underline{\text{Int}}(S_m) \gtrsim \dots \gtrsim \underline{\text{Int}}^n(S_m) \gtrsim \dots \gtrsim \underline{Cv}.$$

If  $S_m$  is not a fixpoint of  $\underline{\text{Int}}$  and the above descending sequence is finite (e.g. the lattice A-Cont satisfies the descending chain condition) its limit is a better approximation of  $\underline{Cv}$  than  $S_m$ . When the sequence is infinite or slowly converging, one can among other solutions approximate its limit.

#### 9.3.1 Decreasing Approximation Sequence

At step  $n$  in the descending sequence, we have :

$$\underline{\text{Int}}^{n-1}(S_m) \gtrsim \underline{\text{Int}}^n(S_m) \gtrsim \underline{Cv}$$

In order to accelerate the convergence, we should for the next step find an approximation  $D$  such that  $\underline{\text{Int}}^{n+1}(S_m) \gtrsim D \gtrsim \underline{Cv}$ . But not knowing  $\underline{Cv}$ , this characterization is very weak since  $D$  could be chosen incorrectly that is to say less than  $\underline{Cv}$  or non comparable with  $\underline{Cv}$ . The fact that  $\underline{Cv}$  is the greatest lower bound of the set of  $X \in A\text{-Cont}$  such that  $\underline{\text{Int}}(X) \lesssim X$  gives a correctness criterion for the choice of  $D$  when  $\underline{Cv}$  is unknown, we must have :

$$\underline{\text{Int}}^{n+1}(S_m) \gtrsim D \gtrsim \underline{\text{Int}}(D)$$

On the contrary to 9.1.1, this characterization does not provide an efficient construction of  $D$ .

#### 9.3.2 Truncated Decreasing Sequence

In front of these difficulties we will enforce convergence by choosing  $D$  such that :

$$\exists n \geq 0 \mid \underline{\text{Int}}(S_m) \gtrsim D \gtrsim \underline{\text{Int}}^{n+1}(S_m)$$

(However, we will not artificially truncate the decreasing sequence by imposing an arbitrary upper bound on  $n$ ).

Let  $\widetilde{\text{D-int}} : \text{A-Cont} \rightarrow \text{A-Cont}$  be such that :

- 9.3.2.1  $\{\forall C \in \text{A-Cont}\} \{C \geq \widetilde{\text{Int}}(C)\} \Rightarrow \{C \geq \widetilde{\text{D-int}}(C) \geq \widetilde{\text{Int}}(C)\}$
- 9.3.2.2  $\forall C \in \text{A-Cont}$ , every infinite sequence  $C, \widetilde{\text{D-int}}(C), \dots, \widetilde{\text{D-int}}^n(C), \dots$  is not strictly decreasing.

The truncated decreasing sequence  $S'_0, \dots, S'_n, \dots$  is recursively defined by :

- 9.3.2.3  $S'_0 = s_m$   
 $S'_{n+1} = \text{if } (S'_n \neq \widetilde{\text{Int}}(S'_n)) \text{ and } (S'_n \neq \widetilde{\text{D-int}}(S'_n))$   
 $\quad \quad \quad \text{then } \widetilde{\text{D-int}}(S'_n)$   
 $\quad \quad \quad \text{else } S'_n$   
 $\quad \quad \quad \text{fi}$

Let us now prove that the truncated decreasing sequence is a finite strictly decreasing chain which terms are greater than  $\underline{\text{Cv}}$  the least fixpoint of  $\widetilde{\text{Int}}$ .

Let  $p$  be the least natural number (eventually infinite) such that  $S'_p = S'_{p+1}$ . Trivially from 9.1.1 :

$$S'_0 = s_m \geq \widetilde{\text{Int}}(S'_0) \geq \underline{\text{Cv}}$$

If  $p > 0$  then  $S'_0 \neq \widetilde{\text{Int}}(S'_0)$ , therefore  $S'_0 \geq \widetilde{\text{D-int}}(S'_0)$ . Then applying 9.3.2.1 we have :

$$S'_0 \geq \widetilde{\text{D-int}}(S'_0) = S'_1 \geq \widetilde{\text{Int}}(S'_0) \geq \underline{\text{Cv}}$$

But 9.3.2.3 implies  $S'_0 \neq \widetilde{\text{D-int}}(S'_0)$ , hence :

$$S'_0 > S'_1 \geq \widetilde{\text{Int}}(S'_0) \geq \underline{\text{Cv}}$$

For the induction step, let us suppose that for  $k < p$ , we have :

$$S'_{k-1} \geq S'_k \geq \widetilde{\text{Int}}(S'_{k-1}) \geq \underline{\text{Cv}}$$

Since  $\widetilde{\text{Int}}$  is order preserving we have :

$$\widetilde{\text{Int}}(S'_{k-1}) \geq \widetilde{\text{Int}}(S'_k) \geq \widetilde{\text{Int}}^2(S'_{k-1}) \geq \widetilde{\text{Int}}(\underline{\text{Cv}}) = \underline{\text{Cv}}$$

By transitivity  $S'_k \geq \widetilde{\text{Int}}(S'_k)$  and since 9.3.2.3 implies  $S'_k \neq \widetilde{\text{Int}}(S'_k)$  we have from 9.3.2.1 :

$$S'_k \geq \widetilde{\text{D-int}}(S'_k) = S'_{k+1} \geq \widetilde{\text{Int}}(S'_k)$$

Since 9.3.2.3 implies  $S'_k \neq \widetilde{\text{D-int}}(S'_k)$  we have :

$$S'_k \geq S'_{k+1} \geq \widetilde{\text{Int}}(S'_k) \geq \underline{\text{Cv}}$$

By recurrence on  $k$  the result is true for  $k \leq p$ . Moreover 9.3.2.2 implies that  $p$  is finite. Q.E.D.

### 9.3.3 Generalization of Kleene's Descending Sequence

When  $\text{A-Cont}$  satisfies the descending chain condition, one can choose  $\widetilde{\text{D-int}}$  to be  $\widetilde{\text{Int}}$ , in which case the final result  $S'_p = \widetilde{\text{Int}}^p(s_m)$  is a fixpoint greater or equal to the least fixpoint  $\underline{\text{Cv}}$  of  $\widetilde{\text{Int}}$ .

The limit of the descending sequence  $S'_0 = \tau, \dots, S'_p = \widetilde{\text{D-int}}^p(\tau), \dots$  is an upper bound of the greatest fixpoint of  $\widetilde{\text{Int}}$ .

### 9.3.4 Narrowing in Truncated Decreasing Sequences

By analogy with 9.1.3 we define now the narrowing operation in order to build a possible construction of  $\widetilde{\text{D-int}}$  by local modifications to  $\widetilde{\text{Int}}$  :

- 9.3.4.1  $\Delta : \text{A-Cont} \times \text{A-Cont} \rightarrow \text{A-Cont}$

- 9.3.4.2  $\forall (C, C') \in \text{A-Cont}^2, \{C \geq C'\} \Rightarrow \{C \geq C \Delta C' \geq C'\}$

- 9.3.4.3 Every infinite sequence  $s_0, \dots, s_n, \dots$  of the form  $s_0 = C_0, s_1 = s_0 \Delta C_1, \dots, s_n = s_{n-1} \Delta C_n, \dots$  for arbitrary abstract contexts  $C_0, C_1, \dots, C_n, \dots$  is not strictly decreasing.

The approximated interpretation

$\widetilde{\text{D-int}} : \text{Arcs}^0 \times \text{A-Cont} \rightarrow \text{A-Cont}$  is defined by :

- 9.3.4.4  $\widetilde{\text{D-int}} = \lambda(q, \underline{\text{Cv}}) . \text{if } q \in \text{W-arcs} \text{ then } \underline{\text{Cv}}(q) \Delta \widetilde{\text{Int}}(q, \underline{\text{Cv}})$   
 $\quad \quad \quad \text{else } \widetilde{\text{Int}}(q, \underline{\text{Cv}})$   
 $\quad \quad \quad \text{fi}$   
 $\text{and } \widetilde{\text{D-int}} = \lambda \underline{\text{Cv}} . (\lambda q . \widetilde{\text{D-int}}(q, \underline{\text{Cv}}))$

This definition of  $\widetilde{\text{D-int}}$  trivially satisfies the requirement 9.3.2.1 since  $\forall C \in \text{A-Cont}$  with property  $\underline{\text{Cv}} \geq \widetilde{\text{Int}}(\underline{\text{Cv}})$  implies  $\underline{\text{Cv}}(q) \geq \widetilde{\text{Int}}(q, \underline{\text{Cv}})$ ,  $\forall q \in \text{Arcs}$ . If  $q \in \text{W-arcs}$  then 9.3.4.2 implies that  $\underline{\text{Cv}}(q) \geq \underline{\text{Cv}}(q) \Delta \widetilde{\text{Int}}(q, \underline{\text{Cv}}) = \widetilde{\text{D-int}}(q, \underline{\text{Cv}}) \geq \widetilde{\text{Int}}(q, \underline{\text{Cv}})$ . Otherwise, if  $q \notin \text{W-arcs}$   $\underline{\text{Cv}}(q) \geq \widetilde{\text{Int}}(q, \underline{\text{Cv}}) = \widetilde{\text{D-int}}(q, \underline{\text{Cv}})$ . Hence  $\underline{\text{Cv}} \geq \widetilde{\text{D-int}}(\underline{\text{Cv}}) \geq \widetilde{\text{Int}}(\underline{\text{Cv}})$ .

The proof of termination (requirement 9.3.2.2) is very similar to the one outlined for  $\widetilde{\text{A-int}}$  in section 9.1.3.

### 9.4 Example : Bounds of Integer Variables

Let us come back to example 9.2. The system of equations was :

- (1)  $C1 = [1, 1]$
- (2)  $C2 = C1 \cup C4$
- (3)  $C3 = C2 \cap [-\infty, 100]$
- (4)  $C4 = C3 + [1, 1]$
- (5)  $C5 = C2 \cap [101, +\infty]$

The ascending approximation sequence led to the approximate solution :

- \*  $C1 = [1, 1]$
- $C2 = [1, +\infty]$
- $C3 = [1, 100]$
- \*  $C4 = [2, 101]$
- $C5 = [101, +\infty]$

Let us define the narrowing  $\Delta$  of intervals by :

- $[ , ]$  is the null element of  $\Delta$ .
- $[i, j] \Delta [k, l] =$   
 $\quad \text{if } i = -\infty \text{ then } k \text{ else } \min(i, k) \text{ fi,}$   
 $\quad \text{if } j = +\infty \text{ then } l \text{ else } \max(j, l) \text{ fi}$

Thus narrowing just discards infinite bounds and makes no improvement on finite bounds, it satisfies the requirements of 9.3.4. According to 9.3.4.4 the system of equations is modified by :

$$(2) \quad C2 = C2 \Delta (C1 \cup C4)$$

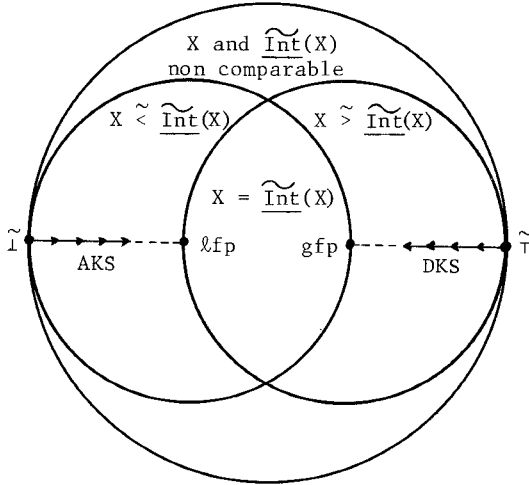
The descending approximation sequence is :

$$\begin{aligned} C2 &= C2 \Delta (C1 \cup C4) \\ &= [1, +\infty] \Delta ([1, 1] \cup [2, 101]) \\ &= [1, +\infty] \Delta [1, 101] \\ * \quad C2 &= [1, 101] \\ C3 &= C2 \cap [-\infty, 100] \\ * \quad C3 &= [1, 101] \cap [-\infty, 100] = [1, 100] \\ &\quad \text{stop on that path.} \\ C5 &= C2 \cap [101, +\infty] \\ * \quad C5 &= [1, 101] \cap [101, +\infty] = [101, 101] \\ &\quad \text{exit.} \end{aligned}$$

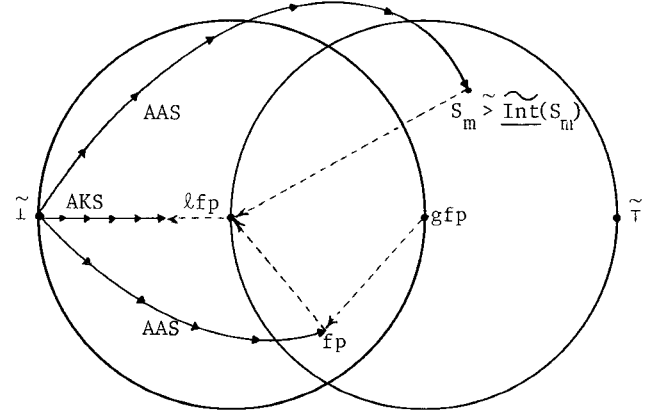
On that example the approximate solution has been improved so that the least fixpoint is reached but this is not the case in general.

### 9.5 Dual Approximation Methods

The lattice  $\widetilde{\text{A-Cont}}$  may be partitioned as follows :

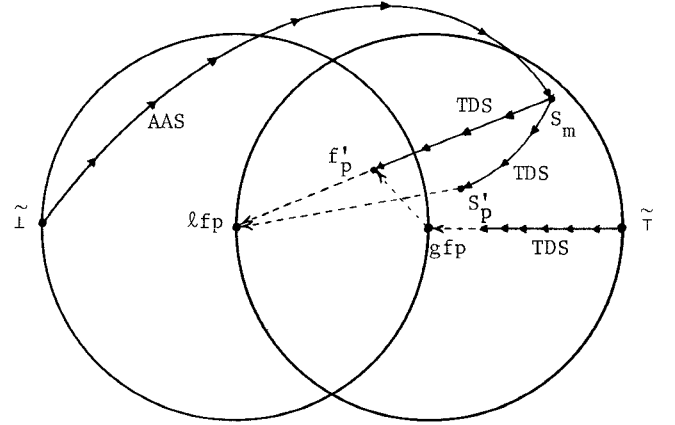


$\text{lfp}$  and  $\text{gfp}$  are the least and greatest fixpoints of  $\widetilde{\text{Int}}$ . The ascending (AKS) and descending (DKS) Kleene's sequences converge toward  $\text{lfp}$  and  $\text{gfp}$  respectively. These limits are reached when  $\widetilde{\text{Int}}$  is continuous. When AKS is infinite we have proposed to use an ascending approximation sequence (AAS) to approximate  $\text{lfp}$ . Its limit may be some fixpoint  $\text{fp}$ , or some  $S_m$  such that  $S_m > \widetilde{\text{Int}}(S_m)$  and  $S_m \gtrsim \text{lfp}$ .

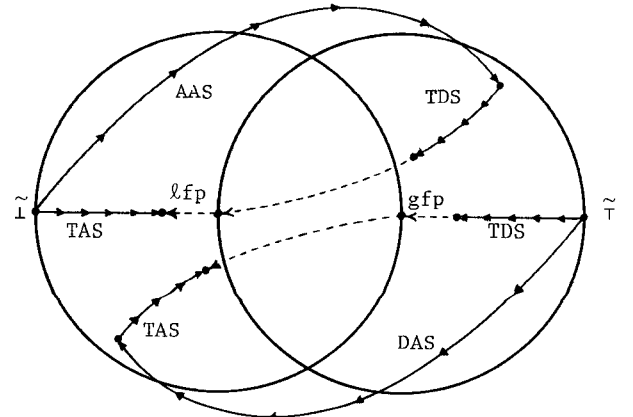


When  $X \gtrsim Y$  we have noted  $X \bullet \text{-----} \bullet Y$ .

The truncated descending sequence TDS is fundamentally different from AAS, since it ensures that the successive approximations starting from  $S_m$  remain in the partition  $\{X \mid X \gtrsim \widetilde{\text{Int}}(X)\}$ , so that their limit  $S'_p$  is greater than  $\text{lfp}$  :



It is clear that the ascending approximation sequence AAS when starting from  $\tilde{1}$  leads to an *upper* bound of the least fixpoint  $\text{lfp}$  of  $\widetilde{\text{Int}}$ , and the truncated descending sequence TDS when starting from  $\tilde{\tau}$  leads to an *upper* bound of the greatest fixpoint  $\text{gfp}$ . Hence the AAS and TDS methods are not dual, therefore when considering their duals DAS and TAS we get a means to surround both extreme fixpoints of  $\widetilde{\text{Int}}$  :



Any of the AAS, TDS, DAS, TAS methods may yields a fixpoint fp which is not the fixpoint lfp or gfp of interest. None of these methods can improve fp to reach lfp or gfp, therefore a "fixpoint improvement method" is necessary. It is our feeling that such a method could be designed only when considering that A-Cont possesses a richer structure (i.e. for particular applications).

Furthermore, in the AAS, TDS, DAS, TAS sequences the term of rank n is computed only as a function of the term of rank n-1, hence these are "separate steps" methods. One can as well imagine to use "bound steps" methods, where the term of rank n is computed as a function of the terms of rank n-1, n-2, ..., n-k. In this last case the Kleene's sequences may be used to compute the first k terms. After k steps more informations about the program would be available to heuristically accelerate the convergence so that the definition of A-int and D-int could be more refined.

Finally, going deeply into the comparism with numerical analysis methods, it is clear that some measure is necessary to control the accuracy of the result. Its definition would certainly also necessitate some additional properties of the abstract contexts.

## 10. Conclusion

It is our feeling that most program analysis techniques may be understood as abstract interpretations of programs. Let us point out global data flow analysis in optimizing compilers (Kildall[73], Morel and Renvoise[76], Schwartz[75], Ullman[75], Wegbreit[75], ...), type verification (Naur[65], ...), type discovery (Cousot[76'], Sintzoff[72], Tenenbaum[74], ...), program testing (Henderson[75], ...) symbolic evaluation of programs (Hewitt et al.[73], Karr[76], ...), program performance analysis (Wegbreit[76], ...), formalization of program semantics (Hoare and Lauer[74], Ligler[75], Manna and Shamir[75], ...), verification of program correctness (Floyd[67], Park[69], Sintzoff[75], ...), discovery of inductive invariants (Katz and Manna[76], ...), proofs of program termination (Sites[74], ...), program transformation (Sintzoff[76], ...), ...

There is a fundamental unity between all these apparently unrelated program analysis techniques : a new interpretation is given to the program text which allows to built an often implicit system of equations. The problem is either to verify that a solution provided by the user is correct, or to discover or approximate such a solution.

The mathematical model we studied in this paper is certainly the weakest which is necessary to unify these techniques, and therefore should be of very general scope. It can be considerably enriched for particular applications so that more powerful results may be obtained.

## Acknowledgements

We wish to thank M. Sintzoff for stimulating discussions. We were very lucky to have F. Blanc do the typing for us.

## 11. References

- Birkhoff[6]. Lattice theory. Amer. Math. Soc. Col. Pub., XXV, Rev. ed.
- Cousot[76]. Static determination of dynamic properties of programs. Programming Symp. Paris. Springer-Verlag Lecture Notes in Comp. Sc. to appear (April).
- Cousot[76']. Static determination of dynamic properties of generalized type unions. Submitted for publication. (Sept.)
- Floyd[67]. Assigning meanings to programs. Proc. Symp. in Appl. Math. Vol. 19. Mathematical Aspects of Computer Science, (J. Schwartz, ed.) AMS, Providence, R.I., 19-32.
- Henderson[75]. Finite state modelling in program development. Proc. Int. Conf. on Reliable Soft., Los Angeles, 221-227.
- Hewitt et al.[73]. Actor induction and meta-evaluation. Conf. Rec. of the First ACM Symp. on Principles of Programming Languages, Boston, 153-168, (Oct.).
- Hoare[67]. An axiomatic basis for computer programming. Comm. ACM 12, 10 (Oct.), 576-580.
- Hoare and Lauer[74]. Consistent and Complementary formal theories of the semantics of programming languages. Acta Inf. 3, 135-153.
- Kam and Ullman[75]. Monotone data flow analysis frameworks. TR.169, C.S. Lab., Princeton Univ.
- Karr[76]. Affine relationships among variables of a program. Acta Inf. 6, 133-151.
- Katz and Manna[76]. Logical analysis of programs. Comm. ACM 19, 4(April), 188-206.
- Kildall[73]. A unified approach to global program optimization. Conf. Rec. of the First ACM Symp. on Principles of Programming Languages, Boston, 194-206, (Oct.).
- Kleene[52]. Introduction to metamathematics. Van Nostrand, New York.
- Ligler[75]. Surface properties of programming language constructs. Int. Symp. on Proving and Improving Programs, (G. Huet and G. Kahn, eds.), IRIA, France.
- Mac Neille[37]. Partially ordered sets. Trans. Amer. Math. Soc., 42, 416-460.
- Manna and Shamir[75]. A new approach to recursive programs. Tech. Rep. CS-75-539, Comp. Sc. Dep., Stanford U.
- Morel and Renvoise[76]. Une méthode globale d'élimination des redondances partielles. Programming Symp. Paris. Springer-Verlag Lecture Notes in Comp. Sc. to appear.(April).
- Naur[65]. Checking of operand types in ALGOL compilers, BIT 5, 151-163.

Park[69]. Fixpoint induction and proofs of program properties. Machine Intelligence 5, (B. Meltzer and D. Michie, eds.), Edinburgh U. Press, 59-78.

Schwartz[75]. Automatic data structure choice in a language of very high level. Comm. ACM 18, 12 (Dec.), 722-728.

Scott[71]. The lattice of flow diagrams. Symp. on Semantics of Programming Languages. Springer-Verlag Lecture Notes in Math. (E. Engeler, ed.), Vol. 188.

Scott and Strachey[71]. Towards a mathematical semantics for computer languages. Tech. Mon. PRG-6, Oxford U. Comp. Lab.

Sintzoff[72]. Calculating properties of programs by valuations on specific models. Proc. ACM conf. on Proving Assertions About Programs. SIGPLAN Notices 7, 1, 203-207.

Sintzoff[75]. Vérifications d'assertions pour les fonctions utilisables comme valeurs affectant des variables extérieures. Int. Symp. on Proving and Improving Programs, (G. Huet and G. Kahn, eds.). IRIA. France.

Sintzoff[76]. Eliminating blind alleys from back-track programs. Proc. of the third Int. Coll. on Automata, Languages and Programming, Edinburgh, (July).

Sites[74]. Proving that computer programs terminate cleanly. Ph.D. Th., Comp. Sc. Dep., Stanford U., (May).

Tarjan[75]. Solving path problems on directed graphs. Tech. Rep. CS-75-528, Comp. Sc. Dep., Stanford U.

Tarjan[76]. Iterative algorithms for global flow analysis. Tech. Rep. CS-76-547, Comp. Sc. Dep., Stanford U.

Tarski[55]. A lattice theoretical fixpoint theorem and its applications. Pacific journal of Math. 5, 285-309.

Tenenbaum[74]. Type determination for very high level languages. NSO-3, Courant Inst. of Math. Sc., New York U., (Oct.).

Ullman[75]. Data flow analysis. Tech. Rep. 179, Dep. of Elec. Eng., Comp. Sc. Lab., Princeton U., (March).

Wegbreit[75]. Property extraction in well-founded property sets. IEEE trans. on Soft. Eng., Vol. SE-1, No. 3, (Sept.).

Wegbreit[76]. Verifying program performance, J. ACM, 23, 4, (Oct.), 691-699.

## 12. Appendix

We note  $\langle L, \cup, \leq, \tau, \perp \rangle$  a complete  $\cup$ -semilattice  $L$ , with partial ordering  $\leq$ , supremum  $\tau$  and infimum  $\perp$ . These definitions are given in Birkhoff[61].

Note :  $L$  is a complete lattice.  
(proof in Birkhoff[61], p. 49).

We take  $f$  is isotone,  $f$  is order-preserving or  $f$  is monotone to be synonymous and mean :

$$\begin{aligned} & \{ \forall (x, y) \in L^2, \{x \leq y\} \implies \{f(x) \leq f(y)\} \} \\ & \iff \{ \forall (x, y) \in L^2, \{f(x \cup y)\} \geq f(x) \cup f(y) \} \end{aligned}$$

(H1) : Let  $F$  be an order-preserving function from the complete semi-lattice  $\langle L, \cup, \leq, \tau, \perp \rangle$  in itself.

( $\overline{H1}$ ) : Let  $\overline{F}$  be an order-preserving function from the complete semi-lattice  $\langle \overline{L}, \cup, \leq, \tau, \perp \rangle$  in itself.

(L1) : The fixpoints of  $F$  form a non-empty complete lattice with supremum  $g$ , infimum  $\ell$  such that :

$$g = \cup \{x \mid (x \in L) \wedge (x \leq F(x))\}$$

$$\ell = \cap \{x \mid (x \in L) \wedge (F(x) \leq x)\}$$

(This result is proved in Tarski[55], pp.286-287). Note that the fixpoints of  $F$  need not form a sublattice of  $L$ .  
We note  $g$  and  $\ell$  the greatest and least fixpoints of  $\overline{F}$ .

(H2) : Let  $\alpha$  and  $\beta$  be such that :

$$(H2.1) \quad \alpha : L \rightarrow \overline{L}$$

$$(H2.2) \quad \gamma : \overline{L} \rightarrow L$$

$$(H2.3) \quad \alpha \text{ is order preserving}$$

$$(H2.4) \quad \gamma \text{ is order preserving}$$

$$(H2.5) \quad \forall \overline{x} \in \overline{L}, \overline{x} = \alpha(\gamma(\overline{x}))$$

$$(H2.6) \quad \forall x \in L, x \leq \gamma(\alpha(x))$$

(H3.1) : (H1), ( $\overline{H1}$ ), (H2) and  $\{\forall x \in L, \overline{F}(\alpha(x)) \geq \alpha(F(x))\}$

(H3.2) : (H1), ( $\overline{H1}$ ), (H2) and  $\{\forall \overline{x} \in \overline{L}, \gamma(\overline{F}(\overline{x})) \geq F(\gamma(\overline{x}))\}$

(L2) :  $\{H3.1\} \iff \{H3.2\}$

Proof :

$$\forall \overline{x} \in \overline{L},$$

$$\overline{F}(\alpha(\gamma(\overline{x}))) \geq \alpha(F(\gamma(\overline{x}))) \text{ by } x = \gamma(\overline{x}) \text{ in H3.1}$$

$$\overline{F}(\overline{x}) \geq \alpha(F(\gamma(\overline{x}))) \text{ from H2.5}$$

$$\gamma(\overline{F}(\overline{x})) \geq \gamma(\alpha(F(\gamma(\overline{x})))) \text{ from H2.4}$$

$$\gamma(\overline{F}(\overline{x})) \geq F(\gamma(\overline{x})) \text{ H2.6 and transitivity.}$$

$$\forall x \in L$$

$$\gamma(\overline{F}(\alpha(x))) \geq F(\gamma(\alpha(x))) \quad \overline{x} = \alpha(x) \text{ in H3.2}$$

$$\gamma(\alpha(x)) \geq x \text{ H2.6}$$

$$F(\gamma(\alpha(x))) \geq F(x) \text{ F order preserving in (H1).}$$

$$\gamma(\overline{F}(\alpha(x))) \geq F(x) \text{ transitivity}$$

$$\alpha(\gamma(\overline{F}(\alpha(x)))) \geq \alpha(F(x)) \text{ H2.3}$$

$$\overline{F}(\alpha(x)) \geq \alpha(F(x)) \text{ H2.5}$$

Q.E.D.

Since H3.1 and H3.2 are proved by L2 to be equivalent, we choose :

(H3) : (H3.1) or (H3.2)

(L3) : Let  $F : L \rightarrow L$  be an order-preserving function from the semilattice  $\langle L, \cup, \leq, \tau, \perp \rangle$  in itself,  $\ell$  and  $g$  respectively the least and greatest fixpoints of  $F$ , then :

$$\forall x \in L, \{g \cup F(x) \geq x\} \iff \{g \geq x\}$$

(The dual of this result is proved in Park[69]. pp. 66). By duality :

$$\forall x \in L, \{\ell \cap F(x) \leq x\} \iff \{\ell \leq x\}$$

(T1) : H1,  $\overline{H1}$ , H2, H3 imply that the greatest fixpoints  $g$  and  $\overline{g}$  of  $F$  and  $\overline{F}$  are related by :

$$\{\alpha(g) \leq \overline{g}\} \text{ and } \{g \leq \gamma(\overline{g})\}$$

Proof :

The existence of  $g$  and  $\overline{g}$  is stated by (L1).

$$\begin{array}{ll} \overline{g} \leq \alpha(g) & \text{trivially} \\ \overline{g} \leq \alpha(F(g)) & \text{since } \overline{g} = F(g) \\ \overline{g} \leq \overline{F}(\alpha(g)) & \text{H3.1, } \overline{U} \text{ isotone, } \geq \text{ transitive} \\ \overline{g} \leq \alpha(g) & \text{L3} \\ \gamma(\overline{g}) \geq \gamma(\alpha(g)) & \text{H2.4} \\ \gamma(\overline{g}) \geq g & \text{H2.6, } \geq \text{ transitive.} \end{array}$$

Q.E.D.

Replacing  $\langle g, \overline{g}, \overline{U}, \geq, \alpha, \gamma, \text{H3.1, H2.4, H2.6} \rangle$  respectively by  $\langle \overline{g}, g, U, \leq, \overline{\alpha}, \overline{\gamma}, \overline{\text{H3.1, H2.4, H2.6}} \rangle$  in the above proof, we get the "dual" theorem :

(T2) : H1,  $\overline{H1}$ , H2, H3 imply that the least fixpoints  $\ell$  and  $\overline{\ell}$  of  $F$  and  $\overline{F}$  are related by :

$$\{\gamma(\overline{\ell}) \geq \ell\} \text{ and } \{\overline{\ell} \geq \alpha(\ell)\}$$

According to Scott[71] a subset  $X \subseteq L$  is called directed if every finite subset of  $X$  has an upper bound (in the sense of  $\leq$ ) belonging to  $X$ . (An obvious example of a directed subset is a non-empty ascending chain). A function  $f : D \rightarrow D$  is called continuous if whenever  $X \subseteq L$  is directed, then  $f(\bigcup\{x \mid x \in X\}) = \bigcup\{f(x) \mid x \in X\}$ .

(H4) : Let  $F$  be a continuous function from the complete semi-lattice  $\langle L, U, \leq, \tau, \perp \rangle$  in itself.

( $\overline{H4}$ ) : Let  $\overline{F}$  be a continuous function from the complete semi-lattice  $\langle \overline{L}, U, \leq, \tau, \perp \rangle$  in itself.

We note  $F^0(x) = x$  and  $F^{n+1}(x) = F(F^n(x))$ .

(L4) :  $H4(\overline{H4})$  implies that  $F(\overline{F})$  has a least fixpoint  $\ell(\overline{\ell})$  which is the limit  $\bigcup_{i=0}^{+\infty} F^i(\perp)$  of the Kleene's sequence  $\perp \leq F(\perp) \leq \dots \leq F^n(\perp) \leq \dots$ .

(The proof is easy to adapt from Kleene[52]'s proof of the first recursion theorem pp. 348-349).

