

Compositional Pointer and Escape Analysis for Java Programs

John Whaley and Martin Rinard

08/16/01

Élodie-Jane Sims. `sims@ens.fr`

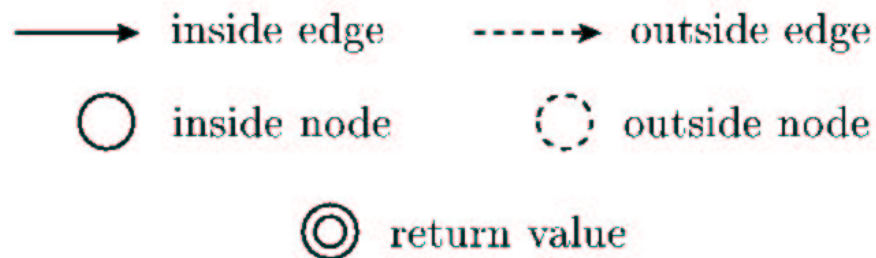
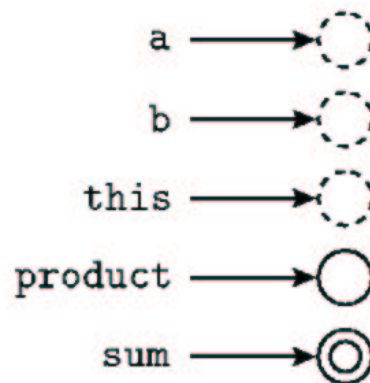
Points-to escape graphs

- Points-to graphs : Characterize how local variables and fields in objects refer to other objects.
- escape information : Characterize how objects allocated in one region of the program can escape to be accessed by an other region.

```

class complex {
  double x,y;
  complex(double a, double b) { x = a; y = b; }
  complex multiply(complex a) {
    complex product =
      new complex(x*a.x - y*a.y,x*a.y + y*a.x);
    return(product);
  }
  complex add(complex a) {
    complex sum = new complex(x+a.x,y+a.y);
    return sum;
  }
  complex multiplyAdd(complex a, complex b) {
    complex product = a.multiply(b);
    complex sum = this.add(product);
    return(sum);
  }
}

```



Properties of the algorithm

- analyze arbitrary part of the program
 - more precise as more of the program is analyzed
 - can distinguish where it does and does not have complete information
- Ô analyze each method independently of its callers
- Ô capable of analyzing a method without analyzing all of the method that it invokes

Applications

- eliminate synchronization for objects that are accessed by only one thread
- allocate objects on the stack instead of in the heap

Escape information

An object can **escape** if it is:

- reachable from a parameter or is the result value of the currently analyzed method
- reachable from a parameter or the result value of an invoked method, and we don't know what the invoked method do with it
- reachable from a static class variable or a runnable object

An object is **captured** if it does not *escape*.

Escape information propagation constraint :

$$\frac{\langle \langle n_1, \mathbf{f} \rangle, n_2 \rangle \in O \cup I}{e(n_1) \subseteq e(n_2)}$$

Program objects

- $v \in V$: variables
 - $l \in L$ local variables
 - $p \in P$ formal parameter variable
 - $cl \in CL$ class names
- $f \in F$: field object (ex: $v.f$ or $cl.f$)
- $op \in OP$: methods
 - cl receiver class
 - $p_0, \dots, p_k : p_0.op(p_1, \dots, p_k)$
- $st \in ST$: nodes of control flow graph,
 $enter_{op}, exit_{op} \in ST$
- $m \in M$: method invocation site
- $n \in N$: nodes

Points-to escaped graph is an abstraction

- nodes represent objects
- edges represent references between objects

Properties:

- ^a a single outside object may be represented by multiple outside nodes
- ^a each object is represented by at most one inside node
- ^a all outside/inside references have a corresponding outside/inside edge in the points-to escape graph
- ^a if an object is represented by a captured node, it is represented by only that node

- a captured objects are reachable only via paths that start with the local variables
- a if a node is captured at the end of a method, the objects that it represents become inaccessible as soon as the method return

The nodes

- N_I : inside nodes
 - represent objects created inside the currently analyzed region and accessed via inside edges
- N_O : outside nodes
 - represent objects created outside the currently analyzed region or accessed via outside edges
 - N_L : load nodes, if $l_1 = l_2.f$ then the nodes pointed by l_1 are load nodes
 - N_P : parameter nodes
 - N_{CL} : class nodes
 - N_R : return nodes, if we have a statement $l_1 = l_2.op(\dots)$ and the analysis skip the call it create a return node to which l_1 will point to.

The edges

- $O \subseteq (N \times F) \times N_L$: outside edges, represent references created outside the currently analyzed method
- $I \subseteq (V \cup (N \times F)) \times N$: inside edges, represent references created inside the currently analyzed method

A points-to escape graph

$$\langle O, I, e, r \rangle$$

u $O \subseteq (N \times F) \times N_L$: outside edges

u $I \subseteq (V \cup (N \times F)) \times N$: inside edges

u $e : N \rightarrow 2^{P \cup CL \cup T \cup M}$: escape function

u $r \subseteq N$: result nodes (\neq *return nodes*)

Points-to escape graph's comments

3 the escape function is just here to avoid to have to compute it when we need escape information, it just follow the definition of *a node can escape*.

3 result nodes \neq return nodes

result nodes = node the analyzed method return (*i.e nodes pointed by l if we have a statement `return l` in the program*)

return nodes = the “wrong” nodes created by skipping call site when analysing this method

Statements

- $l = v$
- $l_1 = l_2.f$
- $l_1.f = l_2$
- `return l`
- $l = \text{new } cl$
- $l = l_0.\text{op}(l_1, \dots, l_k) : \textit{method invocation site}$

Global algorithm

- precompile the program to obtain a control-flow graph of each method with the kind of statement we want
- do a topological sorting of the methods in the order of calling
- analyze the methods in the reverse topological sort order and use a fixed-point algorithm within each strongly connected component
- apply a **dataflow algorithm** on each method to obtain a points-to graph at each program point
- note that the control-flow graph is build without using any information from the test of an `if` or a `while`

Dataflow algorithm for a method

- initialize the points-to information at the entry point of the method
 - the parameters point to the corresponding parameter nodes
 - the class names point to the corresponding class nodes
- **processes the statements** in the method until reaching a fixed point.

Process a statement

- joins the points-to escape graphs flowing into the statement from all of its immediate predecessors in the control-flow graph
- apply the **statement's transfer function** to the joined points-to escape graph

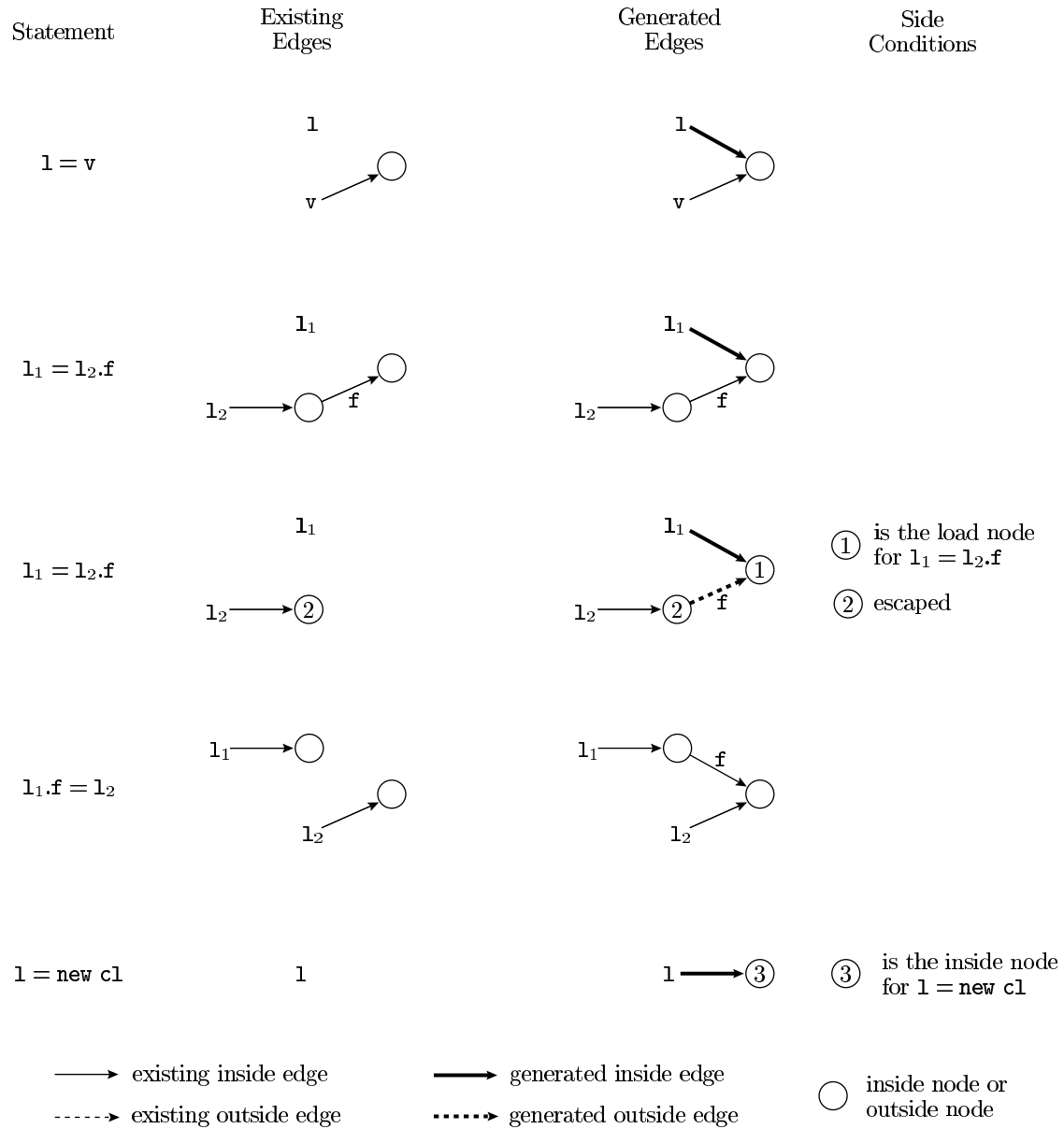
Statement's transfer function

$$I' = (I - \text{Kill}_I) \cup \text{Gen}_I$$

$$O' = O \cup \text{Gen}_O$$

Statement	Kill _I	Gen _I	Kill _O
$l = v$	$\text{edgesFrom}(I, l)$	$\{l\} \times I(v)$	\emptyset
$l_1 = l_2.f, S_E = \emptyset$	$\text{edgesFrom}(I, l)$	$\{l_1\} \times S_I$	\emptyset
$l_1 = l_2.f, S_E \neq \emptyset$	$\text{edgesFrom}(I, l)$	$\{l_1\} \times (S_I \cup \{n\})$	$(S_E \times \{f\}) \times \{n\}$
$l_1.f = l_2$	\emptyset	$(I(l_1) \times \{f\}) \times I(l_2)$	\emptyset
$l = \text{new } cl$	$\text{edgesFrom}(I, l)$	$\{\langle l, n \rangle\}$	\emptyset

- $l_1 = l_2.f : n$ the load node
 $S_E = \{n_2 \in I(l_2).escaped(\langle O, I, e, r \rangle, n_2)\} : \text{set of escaped nodes to which } l_2 \text{ points}$
 $S_I = \cup \{I(n_2.f).n_2 \in I(l_2)\} : \text{set of nodes accessible via inside edges from } l_2.f$
 S
- return $l : r' = I(l)$
- $l = l_0.op(l_1, \dots, l_k) : \text{Call rules}$



Call statement

$$l = l_0.\text{op}(l_1, \dots, l_k)$$

Two choices :

- skipping the call :
 - remove inside edges from l
 - create a new return node for this *method invocation site*
 - add an inside edge from the l to this new return node
 - update the escape function (the parameter escapes by the call site)
- analyzing the call : **mapping process**
 - map some nodes and edges of the called method's points-to graph to corresponding nodes and edges in the current points-to escapes graph
 - using this mapping add nodes and edges to the current points-to escapes graph

Mapping rules

The mapping algorithm take :

- ' the graph of the caller method with the escape function and with the call skipped : $\langle O, I, e, r \rangle$,
 O are the outside edges, I are the inside edges, e is the escape function, r the set of nodes that represent objects that may be returned by the method
- ' the graph of the callee with the escape function :
 $\langle O_R, I_R, e_R, r_R \rangle$
- à and it build a mapping from the nodes of the callee graph to the nodes of the caller, and build a new graph $\langle O_M, I_M, e_M, r_M \rangle$.

- Initialization : [rule 4]
- Parameters : [rule 1]
- Return values + assign l : [rule 8 + 11]
- Add edges using the mapping : [rule 5, 7]
- Complete the mapping : [rule 3, 6, 7]

$$\frac{0 \leq i \leq k, n \in I(1_i)}{n \in \mu(n_{p_i})} \quad (1)$$

$$\frac{c1 \in \mathbb{CL}}{n_{c1} \in \mu(n_{c1})} \quad (2)$$

$$\frac{\langle \langle n_1, f \rangle, n_2 \rangle \in O_R, \langle \langle n_3, f \rangle, n_4 \rangle \in I, \quad n_3 \in \mu(n_1), n_1 \notin N_I}{n_4 \in \mu(n_2)} \quad (3)$$

$$\frac{O \subseteq O_M \quad I - \text{edgesFrom}(I, 1) \subseteq I_M}{e(n) \subseteq e_M(n) \quad r \subseteq r_M} \quad (4)$$

$$\frac{\langle \langle n_1, f \rangle, n_2 \rangle \in I_R}{(\mu(n_1) \times \{f\}) \times \mu(n_2) \subseteq I_M} \quad (5)$$

$$\frac{\langle \langle n_1, f \rangle, n_2 \rangle \in I_R, n \in \mu(n_1), n_2 \in N_I \cup N_R}{n_2 \in \mu(n_2)} \quad (6)$$

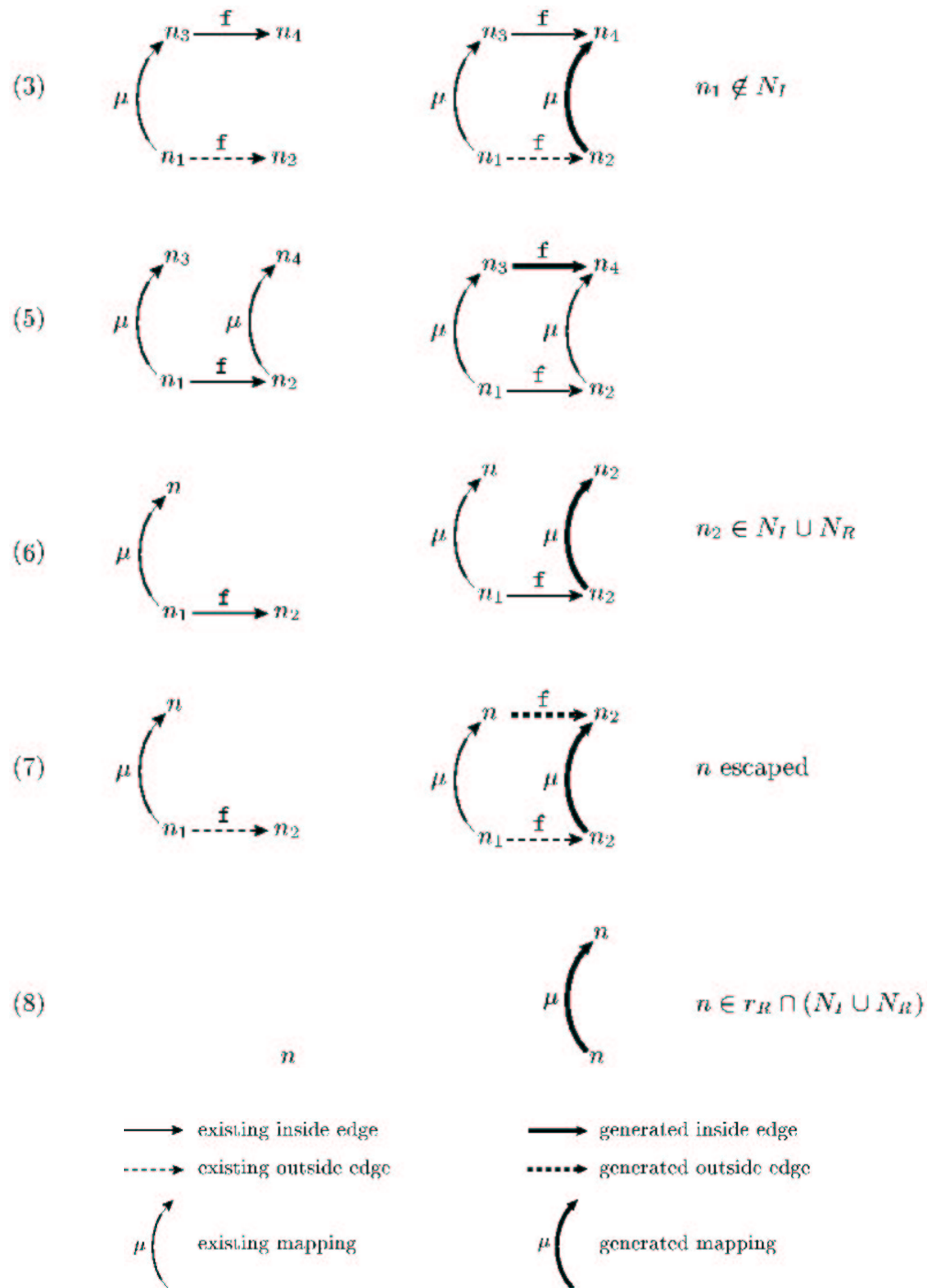
$$\frac{\langle \langle n_1, f \rangle, n_2 \rangle \in O_R, n \in \mu(n_1), \quad \text{escaped}(\langle O_M, I_M, e_M, r_M \rangle, n)}{\langle \langle n, f \rangle, n_2 \rangle \in O_M, n_2 \in \mu(n_2)} \quad (7)$$

$$\frac{n \in r_R \cap (N_I \cup N_R)}{n \in \mu(n)} \quad (8)$$

$$\frac{n' \in \mu(n)}{e_R(n) - P \subseteq e_M(n')} \quad (9)$$

$$\frac{\langle \langle n_1, f \rangle, n_2 \rangle \in I_M \cup O_M}{e_M(n_1) \subseteq e_M(n_2)} \quad (10)$$

$$\cup \{ \mu(n).n \in r_R \} \subseteq I_M(1) \quad (11)$$



Explanation of the rules

- *Initialization* : To build the new graph the algorithm first copy the caller points-to graph but remove from it the insides edges from l (because l will be assigned) [rule 4].
- *Parameters* : Then it links the actual parameters with the formal parameter of the callee [rule 1].
- *Return values + assign l* : We said that the mapping is a kind of equivalence between the nodes of the callee and the nodes of the caller. By the [rule 8], the algorithm map the result nodes of the callee to themselves if they are inside nodes of the callee or if they are return nodes of the callee (*i.e. nodes created by a skipped call when computing the points-to graph of the callee*). Mapping a node to itself “add this node into the new graph”, if we look at the [rule 11], which is used at the end of the algorithm, it add an inside edge from l to all the

nodes “equivalents” of a result node of the callee, so the [rule 8] is a way to make l point to the result nodes which are created by the callee.

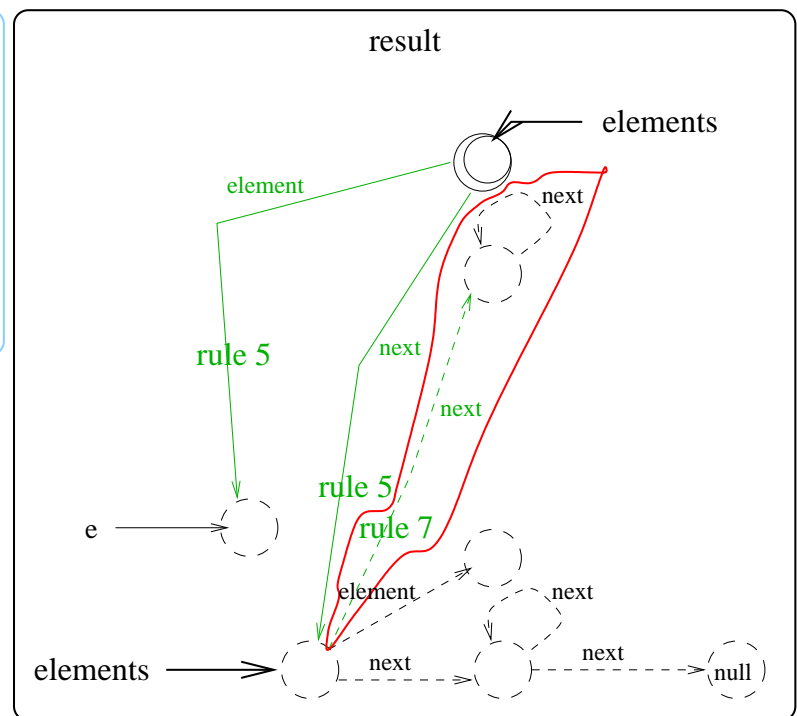
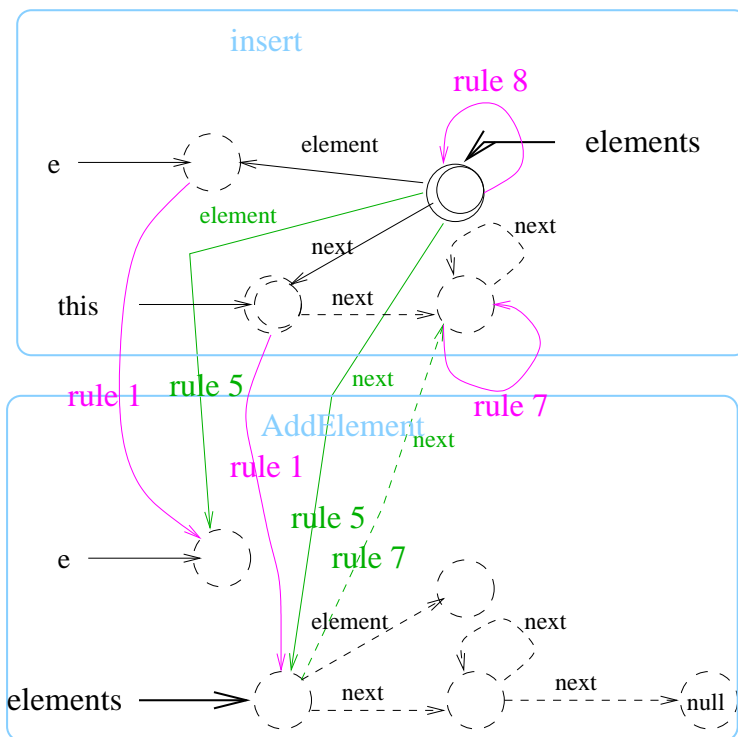
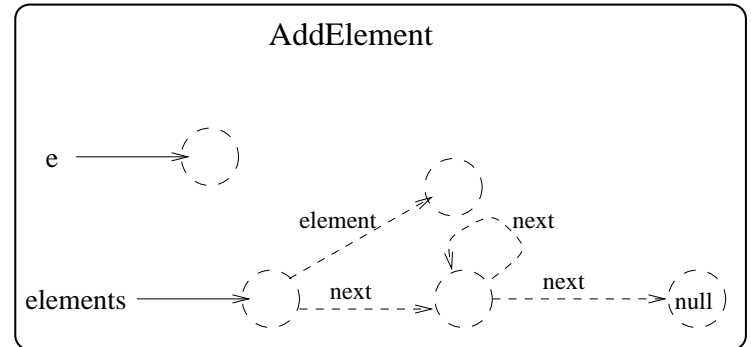
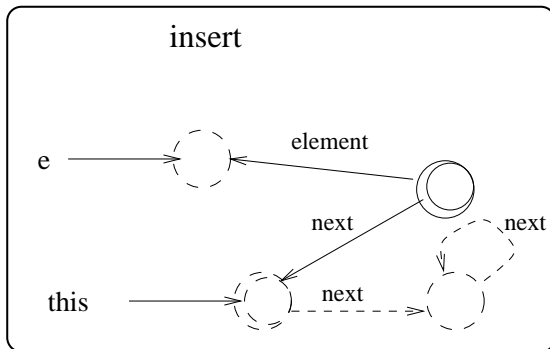
- *Add edges using the mapping :*
 - [rule 5] add an inside field edge to the new graph using the mapping equivalence
 - [rule 7] add an outside edge to the new graph (after having complete the mapping)
- *Complete the mapping :*
 - [rule 3] complete the mapping to add an equivalence between an outside edge of the callee and an inside edge of the caller
 - [rule 6] complete the mapping to add an equivalence between an inside edge of the callee and an inside edge of the caller (*with the help of the [rule 5]*)
 - [rule 7] complete the mapping to add an equivalence between an outside edge of the callee and an outside edge of the caller

Example of application of mapping rules

`elements = elements.insert(e)`

caller : AddElement : l0 = elements, l1 = e

callee : insert : p0 = this, p1=e



**Pb in case $l = l_0.op(l_1, \dots, l_i)$
where there is a k such that $l = l_k$**

```
class list{
list next;
list(){

synchronized list callee (){
next = new list();
return(next);
}
synchronized void caller (){
res = this;
res = res.callee();
}
}
```

