

# Alias analysis

Objetivo: Hacer un análisis del código para modelar aliasing

aliasing: multiples variables apuntan a la misma posición en memoria de forma directa o indirecta. (Para la interpretación de modelado: multiples variables representan la misma locación mutable)

El análisis puede ser may alias (pueden cumplirse para algún camino) o must alias (se cumple para todo camino)

## Como inferir aliasing

### May alias vs Must alias

- May:

- Pares que pueden cumplirse para algún camino
  - (a,b), (a,c), (c,d), (a,d)

```
object a,b,c,d;  
c = d  
if (B)  
  a = b;  
else  
  a = c;
```

- Must:

- Pares que deben cumplir para todo camino
  - (c,d)

- Dataflow:  $Out(N) = Gen(N) + (In(N) - Kill(N))$ ;

- $a = b$ 
  - $Gen = \{ (a,x) \mid (b,x) \in In(N) \}$ ;  $Kill = \{ (a,?) \}$
- $a = new A()$ ;
  - $Gen = \{ \}$ ;  $Kill = \{ (a,?) \}$
- $a = b.f?$   $a.f=b?$

Necesitamos recordar al objeto de tipo A?

- Necesitamos algún modelo de la memoria

Ignoramos los campos?

Con este modelado de generar tuplas no podemos modelar campos ( $a = b.f$ ;  $a.f = b$ ). Para ello necesitamos alguna forma de modelar la memoria (en particular el heap).

## Points-to analysis

Objetivo: armar un análisis que nos diga a donde apunta cada puntero

Un análisis de points-to es un may forward dataflow analysis que permite determinar a que objetos puede apuntar una variable durante la ejecución de un programa.

Se simplifica a 4 operaciones

```
x = new C(); // new  
x = y // copy  
x = y.f // load  
x.f = y // store
```

### Tipos de analisis

- Sensibles a flujo: computan un análisis points-to para cada punto del programa (son caros)
- Insensibles al flujo: hacen un análisis points-to para todo el programa en general y todas las asignaciones son no destructivas. Hay 2 tipos
  - Andersen: basado en inclusión
  - Steengard: basado en unificación (no entra en el parcial)
- Sensibles a contexto: Pueden distinguir entre distintos llamados a metodos.

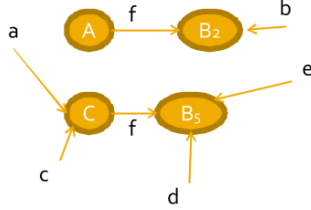
## Points-to Graph (PTG)

Grafo donde tengo nodos por cada objeto (basicamente por cada new), nodos por cada puntero (o variable de tipo objeto, basicamente para trackear el aliasing), y ejes según las relaciones entre objetos, sus campos, y variables (punteros).

Para implementar el análisis utilizaremos un *points-to graph* (PTG) definido como  $G = \langle N, E, L \rangle$  donde

- $N$  es el conjunto de nodos del grafo, y representa los objetos creados por el programa. Cada nodo representa *todos* los objetos creados por una sentencia **new** del programa. Es decir, se creará un solo nodo incluso si el **new** se ejecuta dentro de un ciclo.
- $E$  es el conjunto de aristas del grafo, y representa los objetos alcanzables desde cada objeto mediante campos. Una arista  $e = (n_1, f, n_2)$  en el grafo indica que el nodo  $n_2$  puede ser alcanzado desde el nodo  $n_1$  mediante el campo  $f$ . Tener en cuenta que  $n_1$  y  $n_2$  representan ambos un conjunto de objetos.
- $L : Local \rightarrow P(N)$  es una función de etiquetado que asocia a cada variable local del programa un conjunto de nodos del grafo. Es decir,  $L(x)$  es el conjunto de nodos a los que puede apuntar la variable  $x$ .

Gráficamente, un  $PTG = \langle \{A, B_2, B_5, C\}, \{(A, f, B_2), (C, f, B_5)\}, [a \leftarrow \{C\}, b \leftarrow \{B_2\}, c \leftarrow \{C\}, d \leftarrow \{B_5\}, e \leftarrow \{B_5\}] \rangle$  luce de la siguiente forma:



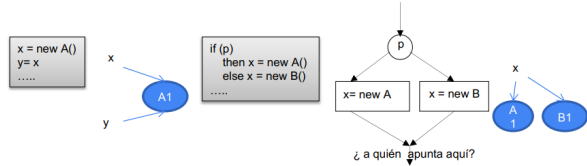
El PTG debe ser utilizado como un reticulado por lo que hay que definir sus operaciones básicas:

- $\perp = (\{\}, \{\}, L)$ : con  $L(x) = \{\}$  para toda variable local.
- $G_1 \subseteq G_2$  si  $L_1 \subseteq L_2$  ó  $(L_1 = L_2 \wedge N_1 \subseteq N_2)$  ó  $(L_1 = L_2 \wedge E_1 \subseteq E_2)$ .
- $G_1 \cup G_2 = \langle N_1 \cup N_2, E_1 \cup E_2, L_1 \cup L_2 \rangle$ .

#### 4ta operación dataflow

El weak update se justifica porque se tiene un análisis may. En el siguiente caso tenemos que  $x$  apunta a  $A$  y  $B$ . Luego si hacemos  $x.f = y$  no sabemos si actualizar el campo  $f$  de  $A$  o  $B$ . Por eso ponemos ambas aristas.

- Si el points-to graph tiene ejes  $(a, A)$  y  $(a, B)$  una variable puede apuntar a  $A$  o  $B$  (MAY)
- Se puede mejorar con sensibilidad a caminos



Las reglas Dataflow para el análisis son las siguientes:

Expresión	Efecto
$p : x = \text{new } A()$	$G' = G \text{ con } L'(x) = \{p\}$ , donde $p$ es el número de línea en el programa.
$x = y$	$G' = G \text{ con } L'(x) = L(y)$
$x = y.f$	$G' = G \text{ con } L'(x) = \{n_2 \mid (n_1, f, n_2) \in E \wedge n_1 \in L(y)\}$
$x.f = y$	$G' = G \text{ con } E' = E \cup \{(n_1, f, n_2) \mid n_1 \in L(x) \wedge n_2 \in L(y)\}$

Pueden asumir que no hay variables globales.

Resúmen de las reglas en formato gráfico:

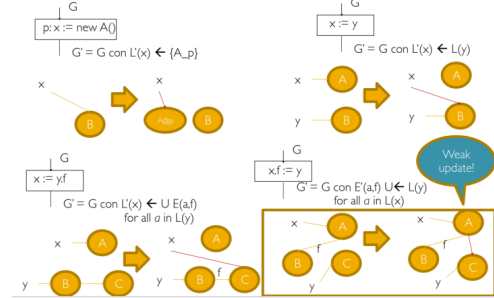
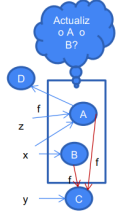


Figura 1: Weak update: No se sabe exactamente qué objeto va a ser actualizado; podría llegar por diferentes caminos.

- Strong update:
  - Cuando se sabe exactamente qué objeto está siendo escrito
  - Ejemplo:  $x := y$
  - Se puede reemplazar la información de points-to de  $x$

- Weak update:
  - No se sabe exactamente qué objeto va a ser actualizado.
  - Ejemplo:  $x.f := y$
  - Durante el análisis puede no saberse a donde podría apuntar  $x$

- Pequeña mejora: si el grado de salida es 1 se puede aplicar un strong update



El PTG como reticulado calcula un PTG por cada punto del programa empezando con el grafo vacío y siguiendo la regla de inclusión definida en las operaciones del reticulado (imagen).

#### Abstracciones del heap

Podemos modelar el heap como un solo nodo, un nodo por cada allocation site (cada new). Este ultimo soluciona el problema de tener un new dentro de un bucle, en tal caso tendríamos que modelar que cada iteración es un objeto nuevo.

```
...
while(i < k) {
  A a = new A();
  b[i] = a;
  ...
}
...
```

## Un nodo por allocation site

```

class L {
  N first;
  A data;
  L(int k) {
    0: N n = new NO;
    1: first = n;
    2: N cur = first;
    3: for(int i=1; i<k; i++) {
    4:   N n2 = new NO;
    5:   cur.next = n2;
    6:   cur = cur.next;
    }
  }
  void m0(L l) {
    0: m1();
    1: l.m1();
  }
  void m1(L l2) {
    0: N cur = first;
    1: while(cur!=null) {
    2:   A a = new AO;
    3:   cur.data = a;
    4:   cur = cur.next;
    }
  }
  void main() {
    0: L l1 = new L(4);
    1: L l2 = new L(1);
    2: l1.m0(l2);
  }
}

```

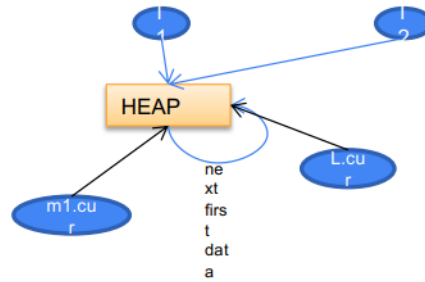
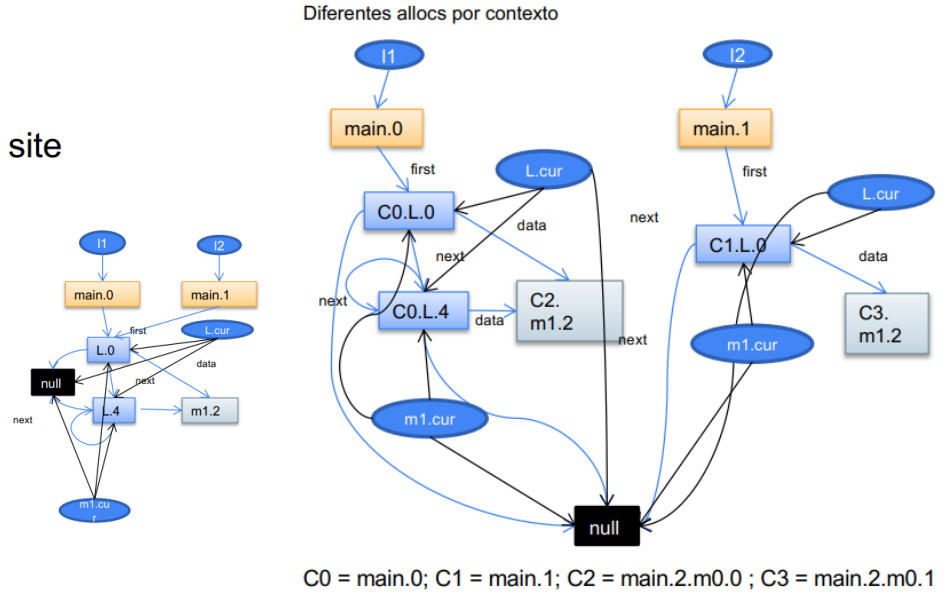


Figure 1: alt text

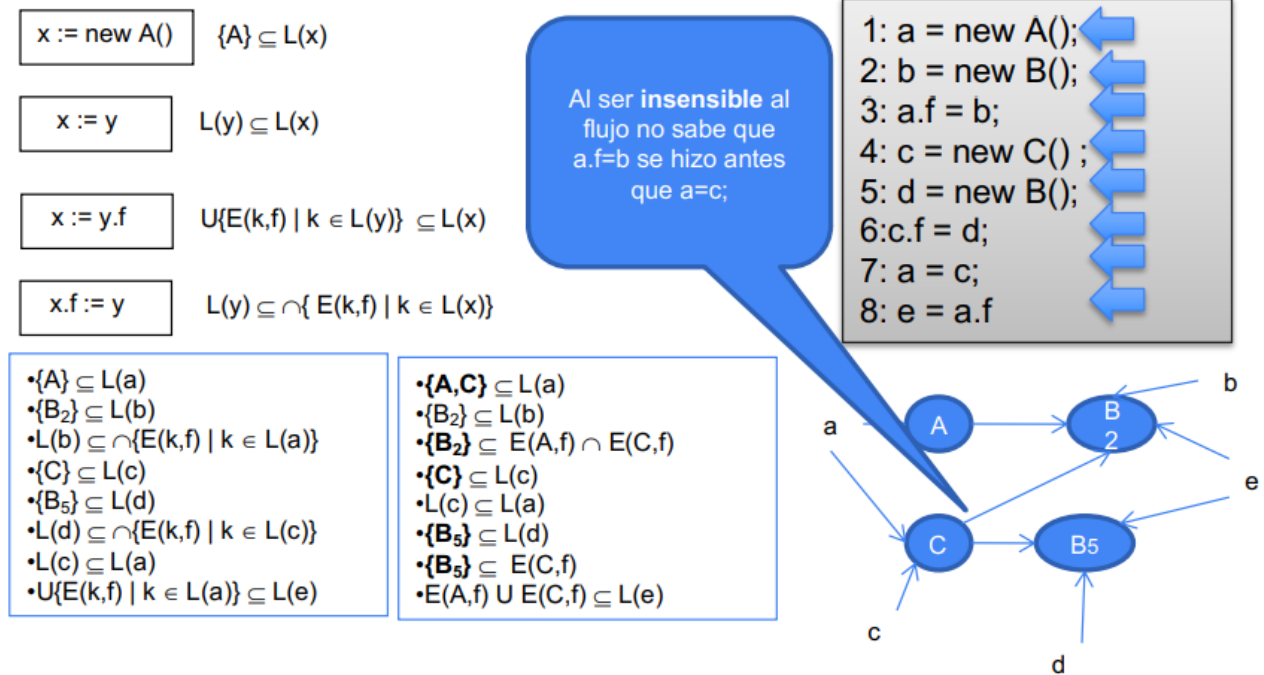
## Andersen

Se basa en inclusiones:  $p = q, pt(q) \subseteq pt(p)$ . Es insensible a contexto y flujo (analiza todo el programa de una). Analisis may.

El algoritmo consiste en recorrer el codigo calculando las restricciones y luego resolverlas con la operación de inclusión hasta llegar a un punto fijo.

$X = \text{alloc } P: \quad \text{alloc-}i \in [X]$   
 $X_1 = \&X_2: \quad X_2 \in [X_1]$   
 $X_1 = X_2: \quad [X_2] \subseteq [X_1]$   
 $X_1 = *X_2: \quad c \in [X_2] \implies [c] \subseteq [X_1] \text{ for each } c \in \text{Cell}$   
 $*X_1 = X_2: \quad c \in [X_1] \implies [X_2] \subseteq [c] \text{ for each } c \in \text{Cell}$

## • Usando inclusiones de restricciones



En palabras, cada instrucción define restricciones a cumplir sobre los objetos del programa:

$x = \text{new}(A)$ : A pertenece a los objetos apuntados por  $x$

$x = y$ : el conjunto de objetos apuntados por  $x$  deben tener a todos los objetos apuntados por  $y$ . (No se borran aristas, en la imagen justo se tapa pero la arista  $a \rightarrow A$  no se borró).

$x = y.f$ : los objetos apuntados por  $x$  deben tener a todos los objetos alcanzables por  $y$  a travez de  $f$

$x.f = y$ : todos los objetos apuntas por  $y$  tienen que pertenecer a la intersección del conjunto de objetos formado por los objetos alcanzables por  $x$  a travez de  $f$