## Pointer and Alias Analysis

**Aliases**:
   two expressions that denote same mutable memory location

Introduced through
- pointers
- call-by-reference
- array indexing
- C unions, Fortran common, equivalence

Applications of alias analysis:
- improved side-effect analysis:
   if assign to one expression,
   what other expressions are modified?
   - if certain modified or not modified, not a problem
   - if uncertain, things can get ugly
- eliminate redundant loads/stores & dead stores
   (CSE & dead assign elim, for pointer ops)
- automatic parallelization of code
   manipulating data structures
- ...

---

## Kinds of alias info

Points-to analysis
- at each program point, calculate set of $p{\rightarrow}x$ bindings,
   if $p$ points to $x$
- two related problems:
   - **may** points-to: $p$ may point to $x$
   - **must** points-to: $p$ must point to $x$

Alias-pair analysis
- at each program point, calculate set of $(\text{expr}_1, \text{expr}_2)$
   pairs, if $\text{expr}_1$ and $\text{expr}_2$ reference the same memory
- **may** and **must** alias-pair versions

Storage shape analysis
- at each program point, calculate an abstract description of
   the structure of pointers etc., e.g. list-like, or tree-like, or
   DAG-like, or ...

Points-to analysis is simple
Alias-pairs analysis more general than points-to analysis,
   but more complicated
Storage shape analysis more abstract

---

## A points-to analysis

At each program point, calculate set of $p{\rightarrow}x$ bindings,
   if $p$ points to $x$

Outline:
- define **may** version first, then consider **must** version
- develop algorithm in increasing stages of complexity
   - pointers only to vars of scalar type
   - add pointers to pointers
   - add pointers to and from structures
   - add pointers to dynamically-allocated storage
   - add pointers to array elements

---

## May-point-to scalars

Domain: Pow(Var $\times$ Var)

Forward flow functions:

$$PT_{p\ :=\ \&x}(in) = in - \{p{\rightarrow}^*\} \cup \{p{\rightarrow}x\}$$

$$PT_{p\ :=\ q}(in) = in - \{p{\rightarrow}^*\} \cup \{p{\rightarrow}v \mid q{\rightarrow}v \in in\}$$

Meet function: union

What about `p := nil`?

## Must-point-to

How to define must-point-to analysis?

Option 1: analogous to may-point-to, but as must problem
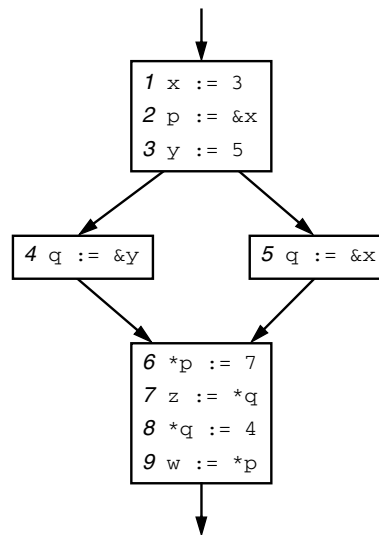- e.g. intersection is meet operation

Option 2: interpretation of may-point-to results
- if $p$ may point to only $x$, then $p$ must point to $x$:

  must-point-to$(p) = \{\, x \mid \{x\} = $ may-point-to$(p)\}$

- what if $p$ may point to $nil$?  $p$ assigned an integer?

## Example

## Using alias info

E.g. reaching definitions

At each program point, calculate set of $x{\to}s$ bindings,
    if $x$ might get its definition from stmt $s$

Simple flow functions:

$RD_{s:\,x\ :=\ \ldots}(\text{in}) = \text{in} - \{x{\to}*\} \cup \{x{\to}s\}$

$RD_{s:\,*p\ :=\ \ldots}(\text{in}) = \text{in} - \{x{\to}* \mid \forall x \in \text{must-point-to}(p)\}$
$\cup \{x{\to}s \mid \forall x \in \text{may-point-to}(p)\}$

## Reaching "right hand sides"

A variation on reaching definitions
    that passes definitions through copies
$x{\to}s$ in set if $x$ might get its definition from rhs of stmt $s$,
    skipping through uninteresting copies and pointer loads
    where possible

Can use reaching right-hand sides to construct def/use chains
    that skip through copies, e.g. for better constant propagation

Additional flow functions:

$RD_{s:\,x\ :=\,y}(\text{in}) = \text{in} - \{x{\to}*\} \cup \{x{\to}s' \mid y{\to}s' \in \text{in}\}$

$RD_{s:\,x\ :=\,*p}(\text{in}) = \text{in} - \{x{\to}*\}$
$\cup \{x{\to}s' \mid p{\to}y \in \text{may-point-to}(p) \land$
$y{\to}s' \in \text{in}\}$

## Another use: "scalar replacement"

If we know that a pointer expression $*p$ aliases a variable $x$
  ($p$ must point to $x$) at some point, then can replace $*p$ with $x$
  • both for load & store

Load part also known as "redundant load elimination"

---

## Adding pointers to pointers

Now allow a pointer to point to a pointer
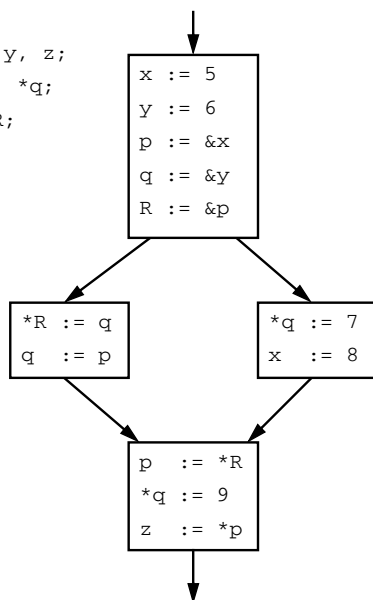  • loads may return pointers, stores may store pointers

New flow functions:
$$PT_{p \ := \ *q}(in) = in - \{p{\to}*\} \cup \{p{\to}v \mid q{\to}r \in in \land r{\to}v \in in\}$$

$$PT_{*p \ := \ q}(in) = in - \{r{\to}* \mid \{r\} = in(p)\}$$
$$\cup \{r{\to}v \mid p{\to}r \in in \land q{\to}v \in in\}$$

---

## Example

```
int x, y, z;
int *p, *q;
int **R;
```



```
x := 5
y := 6
p := &x
q := &y
R := &p
```

```
*R := q
q  := p
```

```
*q := 7
x  := 8
```

```
p  := *R
*q := 9
z  := *p
```

---

## Adding pointers to structs/records/objects/...

A variable can be a structure with a collection of named fields
  • a pointer can point to a field of a structure variable
  • a field can hold a pointer

Introduce location domain: Loc = Var + Loc×Field
  • either a variable or a location followed by a field name
Old PT domain: sets of $v_1{\to}v_2$ pairs = Pow(Var × Var)
New PT domain: sets of $l_1{\to}l_2$ pairs = Pow(Loc × Loc)

Some new forward flow functions:
$$PT_{p \ := \ \&x.f}(in) = in - \{p{\to}*\} \cup \{p{\to}x.f\}$$

$$PT_{p \ := \ x.f} \ (in) = in - \{p{\to}*\} \cup \{p{\to}l \mid x.f{\to}l \in in\}$$
$$PT_{p \ := \ (*q).f}(in) = in - \{p{\to}*\}$$
$$\cup \{p{\to}l \mid q{\to}r \in in \land r.f{\to}l \in in\}$$

$$PT_{x.f \ := \ q} \ (in) = in - \{x.f{\to}*\} \cup \{x.f{\to}l \mid q{\to}l \in in\}$$
$$PT_{(*p).f \ := \ q}(in) = in - \{r.f{\to}* \mid \{r\} = in(p)\}$$
$$\cup \{r.f{\to}l \mid p{\to}r \in in \land q{\to}l \in in\}$$

## Adding pointers to dynamically-allocated memory

```
p := new T
```
- T could be scalar, pointer, structure, ...
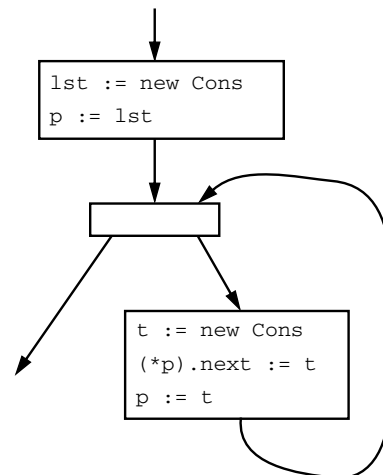
Issue: each execution creates a new location

Idea: generate new var of type T to stand for new location
- make Var domain unbounded
- *newvar*: return next unused element of Var

Flow function:

$PT_{\text{p := new T}}(\text{in}) = \text{in} - \{p \rightarrow *\} \cup \{p \rightarrow newvar\}$

---

## Example



```
lst := new Cons
p := lst
```

```
t := new Cons
(*p).next := t
p := t
```

---

## A monotonic, finite approximation

Can't create a new variable each time analyze statement
- lattice is infinitely tall if Var domain is infinite!
- not a monotonic flow function!

One solution:
    create a special **summary** node for each new stmt
- Loc = Var + **Stmt** + Loc×Field

Fixed flow function:

$PT_{s:\, \text{p := new T}}(\text{in}) = \text{in} - \{p \rightarrow *\} \cup \{p \rightarrow s\}$

Summary nodes abstract a set of possible locations
    $\Rightarrow$ cannot strongly update a summary node

$PT_{*\text{p := q}}(\text{in}) = \text{in} - \{r \rightarrow * \mid \{r\} = \text{in}(p) \land \mathbf{r} \in \mathbf{Loc}\}$
$\qquad\qquad\quad \cup \{r \rightarrow v \mid p \rightarrow r \in \text{in} \land q \rightarrow v \in \text{in}\}$

Alternative summarization strategies:
- summary node for each type T
- *k*-limited summary
  - maintain distinct nodes up to *k* links removed from root vars, then summarize together

---

## Adding pointers to array elements

Array index expressions can generate aliases:
a[i] aliases b[j] if:
- a aliases b and i equals j
- more generally, a and b overlap, and &a[i] = &b[j]

Can have pointers to array elements:
```
p := &a[i]
```

Can have pointer arithmetic, for array addressing:
```
p := &a[0]; ...; p++
```

How to model arrays?

Option 1: reason about array index expressions
    $\Rightarrow$ array dependence analysis

Option 2: use a summary node to stand for all elements of the array