Motivation
A combined abstract domain for lists
Analysis of list-manipulating programs

# Static Analysis of List-Manipulating Programs via Bit-Vectors and Numerical Abstractions

Liqian Chen[1,2]    Renjian Li[1]    Xueguang Wu [1]    Ji Wang[1]

[1]National University of Defense Technology, Changsha, China
[2]National Lab. for Parallel and Distributed Processing, Changsha, China

21/03/2013 – ACM SAC'13 (Track SVT)

Motivation
A combined abstract domain for lists
Analysis of list-manipulating programs

## Overview

- Motivation

- A combined abstract domain for lists

- Analysis of list-manipulating programs

- Conclusion

Motivation
A combined abstract domain for lists
Analysis of list-manipulating programs

# Motivation

Motivation
A combined abstract domain for lists
Analysis of list-manipulating programs

## Motivation

Linked list: a basic dynamic data structure

- commonly used in OS kernels, network protocols, . . .

- **errors**: memory leaks, dangling references, double free, null pointer dereference, . . .

Analysis of list manipulating programs

- **problem**: high complexity

- **solution**: abstraction to make the problem tractable

    - abstraction according to the characteristics of lists
      $\longrightarrow$ simplify the problem & precise enough

    - **shape abstraction + numerical abstraction**
      $\longrightarrow$ numerical related properties over lists

Motivation
A combined abstract domain for lists
Analysis of list-manipulating programs
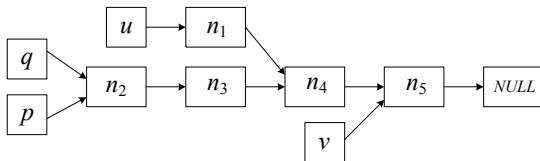
## Motivation

Idea:  combining shape and numerical abstractions
          under the framework of **abstract interpretation**

- a combined abstract domain for singly-linked lists
    - shape: bit vectors
    - numerical: polyhedra, octagons, intervals, . . .

- analysis of list-manipulating programs based on this domain

Motivation
A combined abstract domain for lists
Analysis of list-manipulating programs

shape abstraction for lists
numerical abstraction for lists

# A combined abstract domain for lists

Motivation
A combined abstract domain for lists
Analysis of list-manipulating programs

shape abstraction for lists
numerical abstraction for lists

## Concrete heap state



Shape graph: $\langle N, V, E \rangle$

- $N = \{u, v, p, q, n_1, n_2, n_3, n_4, n_5\}$
- $V = \{u, v, p, q\}$
- $E = \{\langle u, n_1 \rangle, \langle p, n_2 \rangle, \langle q, n_2 \rangle, \langle v, n_5 \rangle, \langle n_1, n_4 \rangle, \langle n_2, n_3 \rangle,$
  $\langle n_3, n_4 \rangle, \langle n_4, n_5 \rangle, \langle n_5, NULL \rangle\}$

Limitations of shape graphs:

- high memory costs (explicit storage)
- lists with symbolic length

Motivation
A combined abstract domain for lists
Analysis of list-manipulating programs

shape abstraction for lists
numerical abstraction for lists

## Shape abstraction for lists
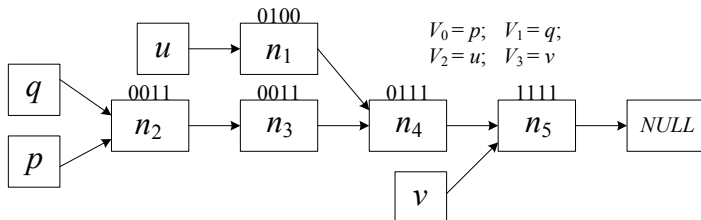
### Definition (*Reach* predicate)

- $Reach(n, n') \triangleq \exists k \in \mathbb{N}. \forall 0 \leqslant i \leqslant k. n_i \in N.$
  $$n_0 = n \wedge n_k = n' \wedge \forall 0 \leqslant j < k. \langle n_j, n_{j+1} \rangle \in E$$
- i.e., $Reach(n, n') = true$ iff there exists a path from $n$ to $n'$

### Definition (Variable Reachability Vector)

For each node $n \in (N - V)$, we define a *Variable Reachability Vector (VRV)* $\mathbf{vec}_n \in \{0, 1\}^{|V|}$ that is a bit-vector of length $|V|$, where

$$\mathbf{vec}_n[i] = 1 \quad iff \quad Reach(V_i, n) = true$$

Motivation
A combined abstract domain for lists
Analysis of list-manipulating programs

shape abstraction for lists
numerical abstraction for lists

## Shape abstraction for lists



Variable Reachability Vector: describe reachability properties of all variables to nodes

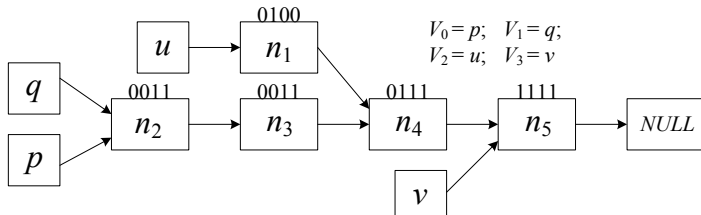- each VRV can be considered as an abstract node

Motivation
A combined abstract domain for lists
Analysis of list-manipulating programs

shape abstraction for lists
numerical abstraction for lists

## Shape abstraction for lists

Reachability information from VRVs

### Definition

- Let $\Gamma$ denote the set of VRVs of all nodes. For arbitrary $\textbf{vec} \in \Gamma$, let $\mathcal{I}_{\textbf{vec}}$ denote the set of the 1-bits in $\textbf{vec}$:
$$\mathcal{I}_{\textbf{vec}} \triangleq \{i \in \mathbb{N} \mid \textbf{vec}[i] = 1\}$$

- If $i \in \mathcal{I}_{\textbf{vec}}$, then $V_i$ can reach (the corresponding nodes) of $\textbf{vec}$, denoted as $V_i \in \textbf{vec}$

E.g., $\mathcal{I}_{0100} = \{2\}; \mathcal{I}_{0011} = \{0, 1\};$

Motivation
A combined abstract domain for lists
Analysis of list-manipulating programs

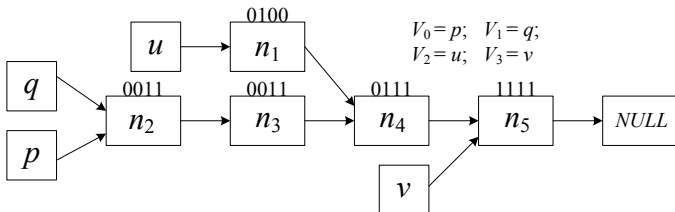shape abstraction for lists
numerical abstraction for lists

# Shape abstraction for lists

## Definition (Reachability relations between VRVs)

Given two VRVs $\mathbf{vec}_1, \mathbf{vec}_2$,

- if $\mathcal{I}_{\mathbf{vec}_1} \subseteq \mathcal{I}_{\mathbf{vec}_2}$, then $\mathbf{vec}_1$ can reach $\mathbf{vec}_2$ ($\mathbf{vec}_1 \subseteq \mathbf{vec}_2$)
- if $\mathcal{I}_{\mathbf{vec}_1} \subset \mathcal{I}_{\mathbf{vec}_2}$, then $\mathbf{vec}_1$ can strictly reach $\mathbf{vec}_2$ ($\mathbf{vec}_1 \subset \mathbf{vec}_2$)
- if $\mathcal{I}_{\mathbf{vec}_1} \cap \mathcal{I}_{\mathbf{vec}_2} = \emptyset$, then $\mathbf{vec}_1$ and $\mathbf{vec}_2$ can not reach each other ($\mathbf{vec}_1 \cap \mathbf{vec}_2 = \emptyset$)
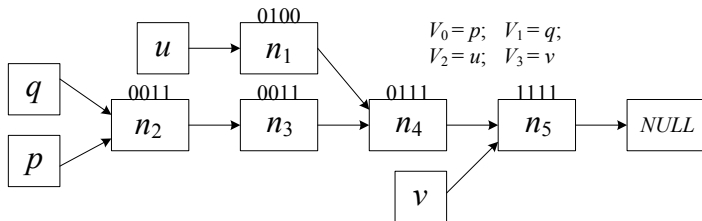
E.g., $\mathbf{vec}_{0100} \subset \mathbf{vec}_{0111}$; $\mathbf{vec}_{0011} \subset \mathbf{vec}_{0111}$; $\mathbf{vec}_{0100} \cap \mathbf{vec}_{0011} = \emptyset$;

Motivation
A combined abstract domain for lists
Analysis of list-manipulating programs

shape abstraction for lists
numerical abstraction for lists
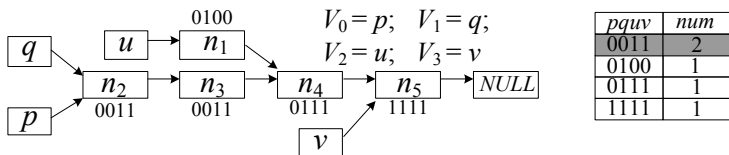
## Properties of VRVs

E.g., Given a set of VRVs $\{0011, 0100, 0111, 1111\}$

- $\mathbf{vec}_{0011} \subset \mathbf{vec}_{0111} \subset \mathbf{vec}_{1111}$
- $p$ points to 0011
- $p, q$ are alias
- $p$ cannot reach the node pointed to by $u$



$V_0 = p; \quad V_1 = q;$
$V_2 = u; \quad V_3 = v$

The set $\Gamma$ of VRVs of a singly-linked list satisfies $|\Gamma| \leq 2|V|$

Motivation
A combined abstract domain for lists
Analysis of list-manipulating programs

shape abstraction for lists
numerical abstraction for lists

## VRVs with counters



$V_0 = p; \quad V_1 = q;$
$V_2 = u; \quad V_3 = v$

| pquv | num |
|------|-----|
| 0011 | 2   |
| 0100 | 1   |
| 0111 | 1   |
| 1111 | 1   |

### Definition

The set of VRVs with counters *VRVCs* $\Gamma^+ \subseteq \Gamma \times \mathbb{N}$ is defined as a set of 2-tuples $\langle \textbf{vec}, num \rangle$ where

- **vec** $\in \Gamma$

- $num \in \mathbb{N}$: the number of the list nodes whose VRV is **vec**

Lists:

- **shape**: VRVs $\leftarrow$ nodes; VRV reachability relations $\leftarrow$ edges

- **numerical**: counters $\leftarrow$ quantitative information of the nodes

Motivation
A combined abstract domain for lists
Analysis of list-manipulating programs

shape abstraction for lists
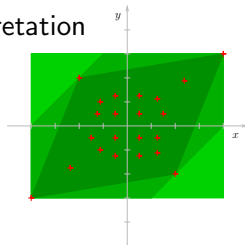numerical abstraction for lists

## Numerical abstraction for lists

Counter variables: auxiliary (non-negative) integer variables

- for each **vec** $\in$ *VRVs*, introduce a counter variable $t^{\textbf{vec}} \in \mathbb{N}$
  - to record the number of the list nodes whose VRV is **vec**
- a special auxiliary variable $t^{0...00} \in \mathbb{N}$
  - to specify memory leak when $t^{0...00} > 0$
- variable ordering: $t^{0...00} \prec t^{0...01} \prec t^{0...10} \prec \cdots \prec t^{1...11}$
- a bijection between **vec** and $t^{\textbf{vec}}$
- $\{\langle \textbf{vec}, t^{\textbf{vec}} \rangle \mid t^{\textbf{vec}} > 0\}$ represents a list, if it is consistent

Motivation
A combined abstract domain for lists
Analysis of list-manipulating programs

shape abstraction for lists
numerical abstraction for lists

# Numerical abstraction for lists

<u>Numerical abstract domains</u> in abstract interpretation

- infer relations among numerical variables
- examples
  - intervals ($a \leq x \leq b$)
  - octagons ($\pm x \pm y \leq c$)
  - polyhedra ($\Sigma_k a_k x_k \leq b$)



<u>Chosen numerical abstract domains</u> for counter variables $t^{\mathbf{vec}}$

- intervals ($a \leq x \leq b$)
- affine equalities ($\Sigma_k a_k x_k = b$)

Motivation
A combined abstract domain for lists
Analysis of list-manipulating programs

Example analysis

# Analysis of list-manipulating programs

Motivation
A combined abstract domain for lists
**Analysis of list-manipulating programs**
Example analysis

# Analysis of list-manipulating programs

$$
\begin{array}{rcl}
p, q \in PVar & & \\
AsgnStmnt & := & p := null \mid p := q \mid p := q \rightarrow next \mid p \rightarrow next := null \mid \\
& & p \rightarrow next := q \mid p := \mathrm{malloc}() \mid \mathrm{free}(p) \\
Cond & := & p == q \mid p == null \mid \neg Cond \mid Cond_1 \lor Cond_2 \mid \\
& & Cond_1 \land Cond_2 \mid \mathrm{true} \mid \mathrm{false} \mid \mathrm{brandom} \\
BranchStmnt & := & \textbf{if } Cond \textbf{ then } \{Stmnt; \}^* \, [\textbf{else } \{Stmnt; \}^* \,] \textbf{ fi} \\
WhileStmnt & := & \textbf{while } Cond \textbf{ do } \{Stmnt; \}^* \textbf{ od} \\
Stmnt & := & AsgnStmnt \mid BranchStmnt \mid WhileStmnt \\
Program & := & \{Stmnt; \}^*
\end{array}
$$

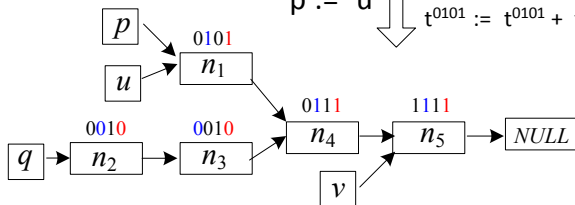Domain operations: on top of shape and numerical abstraction

- inclusion test $\sqsubseteq$
- join $\sqcup$
- widening $\nabla$
- transfer functions $\tau$
  - condition test
  - assignment
- . . .

Motivation
A combined abstract domain for lists
**Analysis of list-manipulating programs**

Example analysis

# $[\![ p := q ]\!]^{\sharp}$

Let $\textbf{vec}' \overset{\text{def}}{=} \textbf{vec}/_{\{p\} \leftarrow q}$. For each $\textbf{vec} \in \Gamma$ s.t. $\textbf{vec}' \neq \textbf{vec}$, we build numerical statements:

$$if(t^{\textbf{vec}} \geq 1)\{ \ t^{\textbf{vec}'} := t^{\textbf{vec}'} + t^{\textbf{vec}}; \ t^{\textbf{vec}} := 0; \ \}$$

Motivation
A combined abstract domain for lists
Analysis of list-manipulating programs

Example analysis

## Example analysis

void copy_and_delete(List* xList) { /* $\heartsuit : \forall t^{vec}.NoOccurenceOf\ vec\ implies\ t^{vec} = 0$ */
  /* assume \length(xList)==9; */
1:   List* yList, pList, qList;       /* $pList \prec qList \prec xList \prec yList$ */
    $/ * t^{0100} = 9; t^{0100} \in [9, 9]; \heartsuit * /$
2:   yList = xList;     qList = pList = null;
    $/ * t^{0100} + t^{1100} = 9, t^{0011} + t^{1100} = 9; t^{0100} \in [1, 9], t^{1100} \in [0, 9], t^{0011} \in [0, 9]; \heartsuit * /$
3:   while (yList != null){
    $/ * t^{0100} + t^{1100} = 9, t^{0011} + t^{1100} = 9; t^{0100} \in [1, 8], t^{1100} \in [1, 9], t^{0011} \in [0, 8]; \heartsuit * /$
4:      yList = yList $\rightarrow$ next; qList = malloc();
5:      qList $\rightarrow$ next = pList; pList = qList;}
    $/ * t^{0011} = 9, t^{0100} = 9; t^{0011} \in [9, 9], t^{0100} \in [9, 9]; \heartsuit * /$
6 :  yList = xList;
    $/ * t^{0011} - t^{1100} = 0; t^{0011} \in [0, 9], t^{1100} \in [0, 9]; \heartsuit * /$
7 :  while (yList != null){
    $/ * t^{0011} - t^{1100} = 0; t^{0011} \in [1, 9], t^{1100} \in [1, 9]; \heartsuit * /$
8:      yList = yList $\rightarrow$ next; qList = qList $\rightarrow$ next;
9:      free(xList); free(pList); xList = yList; pList = qList;
10:  } $/ * \forall \mathbf{vec}.t^{\mathbf{vec}} = 0 * /$
}

Motivation
A combined abstract domain for lists
Analysis of list-manipulating programs

Example analysis

# Example analysis

```
void copy_and_delete(List* xList) { /* ♡ : ∀tᵛᵉᶜ.NoOccurenceOf vec implies tᵛᵉᶜ = 0 */
     /* assume \length(xList)==9; */
1:   List* yList, pList, qList;          /* pList ≺ qList ≺ xList ≺ yList */
2:   yList = xList;    qList = pList = null;
3:   while (yList != null){
4:       yList = yList → next; qList = malloc();
5:       qList → next = pList; pList = qList;}
6 :  yList = xList;
7 :  while (yList != null){
          /*  {t⁰⁰¹¹ − t¹¹⁰⁰ = 0; t⁰⁰¹¹ ∈ [1, 9], t¹¹⁰⁰ ∈ [1, 9]}; ♡ */
8:       yList = yList → next; qList = qList → next;
9:       free(xList); free(pList); xList = yList; pList = qList;
10:  }
}
```

- *pList, qList* are alias; *xList, yList* are alias
- the length of *pList* equals to that of *xList*
- no null pointer dereference

Motivation
A combined abstract domain for lists
Analysis of list-manipulating programs

Example analysis

# Example analysis

void copy_and_delete(List* xList) { /* ♡ : $\forall t^{vec}.NoOccurenceOf\ vec\ implies\ t^{vec} = 0$ */
    /* assume \length(xList)==9; */
1:   List* yList, pList, qList;      /* $pList \prec qList \prec xList \prec yList$ */
2:   yList = xList;    qList = pList = null;
3:   while (yList != null){
4:      yList = yList → next; qList = malloc();
5:      qList → next = pList; pList = qList;
6 :  yList = xList;
7 :  while (yList != null){
8:     yList = yList → next; qList = qList → next;
9:     free(xList); free(pList); xList = yList; pList = qList;
10:  } / * $\forall vec.t^{vec} = 0$; ♡ * /
}

- all heap cells are freed

Motivation
A combined abstract domain for lists
Analysis of list-manipulating programs

Example analysis

# Example analysis

void copy_and_delete(List* xList) { /* ♡ : $\forall t^{vec}.NoOccurenceOf\ vec$ implies $t^{vec} = 0$ */
   /* assume \length(xList)==9; */
1:   List* yList, pList, qList;              /* $pList \prec qList \prec xList \prec yList$ */
2:   yList = xList;    qList = pList = null;
3:   while (yList != null){
4:       yList = yList → next; qList = malloc();
5:       qList → next = pList; pList = qList;
6 :   yList = xList;
7 :   while (yList != null){
8:       yList = yList → next; qList = qList → next;
9:       free(xList); free(pList); xList = yList; pList = qList;
10:   }
}

- Global invariants: $t^{0\cdots0} \equiv 0 \rightsquigarrow$ no memory leak

Motivation
A combined abstract domain for lists
Analysis of list-manipulating programs

Example analysis

# Conclusion

Summary: analysis of lists via abstract interpretation

- **main idea**: combining shape and numerical abstractions
- a combined abstract domain for lists
  - the **structural** information of the shape: bit vectors
    - each bit-vector represents a list segment
  - the **number** of nodes in a segment: numerical abstract domains
    - a counter variable to record the number of nodes in a list segment

Future work

- reasoning over the content of lists (e.g., lists of integers)
  - enable infering advanced properties such as sortedness, no duplicated elements