

# Learning to Decipher the Heap for Program Verification



Marc Brockschmidt<sup>†</sup>, Yuxin Chen<sup>‡</sup>, Byron Cook<sup>†</sup>, Pushmeet Kohli<sup>†</sup>, Daniel Tarlow<sup>†</sup>

<sup>‡</sup> **ETH** zürich

<sup>†</sup> **UCL**

## Motivating Application

```

procedure concat(a: Node, b: Node)
  returns (res: Node)
  requires lseg(a, null) && lseg(b, null)
  ensures lseg(res, null)
{
  if (a == null) {
    return b;
  } else {
    var cur := a;
    while (cur.next != null) {
      cur := cur.next;
    }
    cur.next := b;
    return a;
  }
}
    
```

**Execute**

**Memory safe?**

**GRASShopper** **Instrument**  
**program**

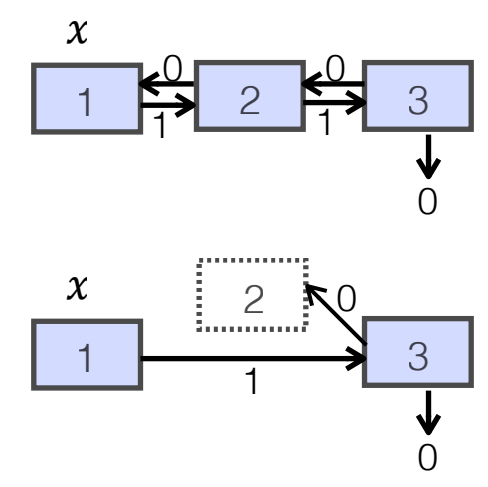
$ls(cur, null) * ls(a, cur) * ls(b, null)$

**Predict**

```

// Remove an element from a doubly-linked list
void RemoveItemFromDLL(DLL* x, int itm)
{
  // if item is at 1st node
  if (x->item == itm) { /* ... */ }
  node *current, *del = x;
  while (del->next != NULL && del->item != itm) {
    del = del->next;
  }
  // if target not found || target at the end
  if (del->next == NULL) { /* ... */ }
  current = del->previous;
  current->next = del->next;
  del->next->previous = current;
  delete(del);
}
    
```

An example code snippet  
which is **memory unsafe**



## Problem Description

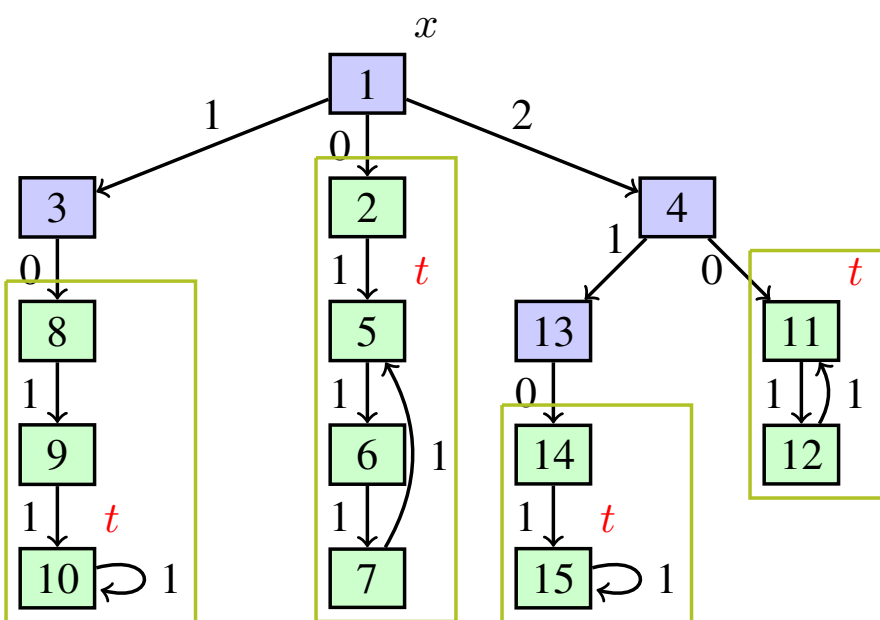
To find **formal descriptions** of the **data structures** that are instantiated, which are used as input to a proof procedure that verifies the program.

## Input Representation

### Heap graphs

**Nodes:**  
address in memory at which **a sequence of pointers** are stored

**Edges:**  
the values of these pointers



Binary tree of "panhandle lists"

## Output Representation

**logical description** of the instantiated data structures, (i.e., formal description of a set of allowed heaps)

**Inductive predicates:**

e.g.,  $ls(x, y, \varphi)$ ,  $tree(x, \varphi)$

Binary tree of "panhandle lists"

$\psi = tree(x, \lambda i_1, i_2, i_3, i_4 \rightarrow$   
 $\exists t. ls(i_2, t, \lambda i_5, i_6, i_7, i_8 \rightarrow \top)$   
 $* ls(t, t, \lambda i_9, i_{10}, i_{11}, i_{12} \rightarrow \top)).$

Given a predicate, we allow  
**nested subformulas**

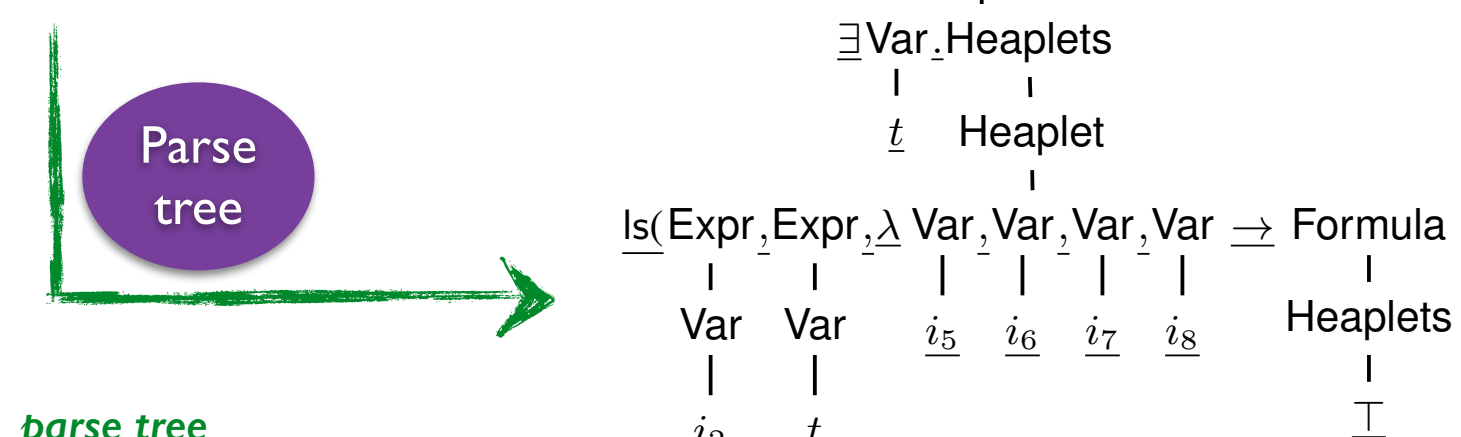
## Problem Formalization

To automatically predict a separation logic formula from a given heap  $H$

### Separation logic formulas

Formula  $\rightarrow$  Heaplets  $| \exists \text{Var. Heaplets} | \exists \text{Var. Var. Heaplets} | \dots$   
Heaplets  $\rightarrow \top | \text{Heaplet} * \text{Heaplets}$   
Heaplet  $\rightarrow ls(\text{Expr}, \text{Expr}, \lambda \text{Var. Var. Var. Var} \rightarrow \text{Formula})$   
 $| tree(\text{Expr}, \lambda \text{Var. Var. Var. Var} \rightarrow \text{Formula})$   
Expr  $\rightarrow \text{NULL} | \text{Var}$

$\exists t. ls(i_2, t, \lambda i_5, i_6, i_7, i_8 \rightarrow \top)$



parse tree

$\mathcal{T} = (\mathcal{A}, g(\cdot), ch(\cdot))$

$\mathcal{A} = \{1, \dots, A\}$

$nodes \mapsto \text{terminal / non-terminal}$   
in separation logic grammar  
 $g: \mathcal{A} \mapsto \mathcal{S}$

$nodes \mapsto \text{children tuple}$   
 $ch: \mathcal{A} \mapsto \mathcal{A}^*$

Partial tree  $\mathcal{T}_{<a}: \text{Parse tree restricted to nodes } \{1, \dots, a\}$

## Problem Statement

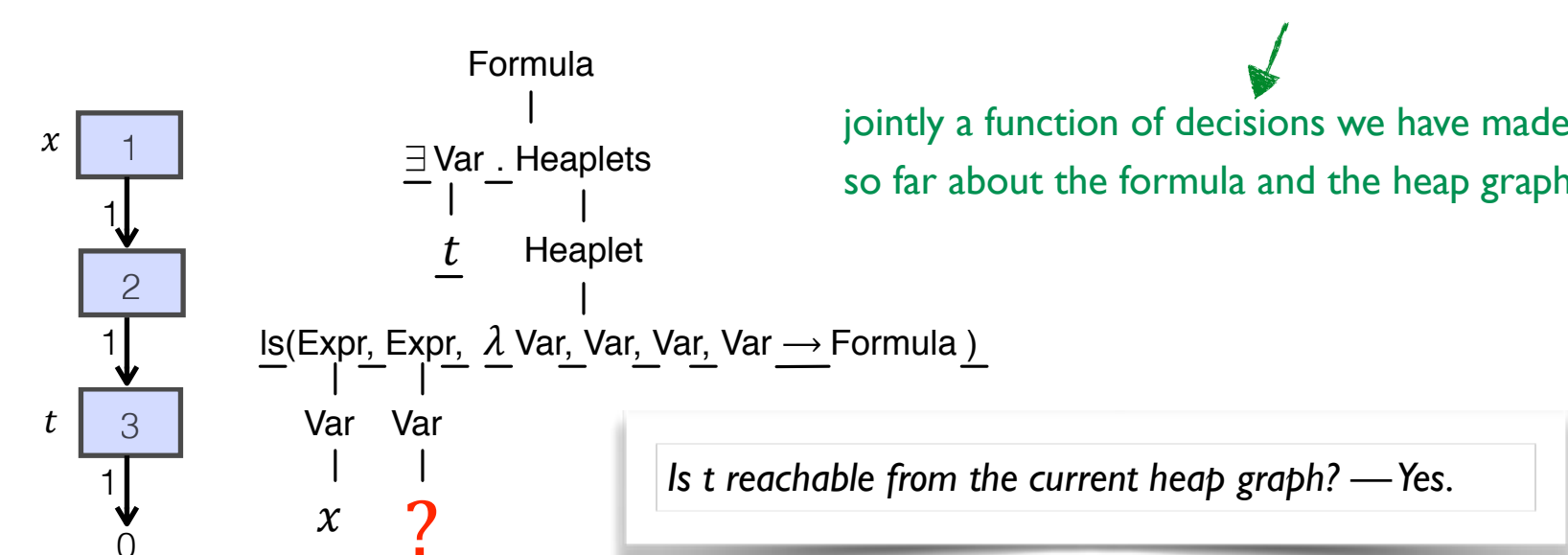
Sequentially predict the next node, conditional upon everything that has been predicted so far:

$$P(\mathcal{T}) = \prod_{a: g(a) \in \mathcal{N}} p(ch(a) \mid H, \mathcal{T} < a).$$

Maximal likelihood training  
 $\Downarrow$   
Independent classification  
Nonterminal  $\mapsto$  Children tuple

## Features

Heap features	Partial tree features	Joint features
(int list) <b>l-gram</b> (frequency of different pointer types)	(int) <b>depth</b> depth of the current nonterminal node	(bool) <b>if variable "y" is reachable</b>
(int list) <b>2-gram</b> "adjacent" pointers	(int) <b>frequency</b>	(bool) <b>if variable "y" is accounted for</b>
(bool) <b>graphic is cyclic</b>	(string) <b>predecessor token</b>	(string) <b>relation to defined predicates</b>



jointly a function of decisions we have made so far about the formula and the heap graph

## Experimental Results

**Prediction model:** Random Forests, **Multiclass Logistic Regression**, Multiclass Neural Network

**Predicates used:** lists, cyclic lists, panhandle lists and trees.

Nesting level	Number of free variables	Number of Formulas
1	1	127
1	2	33254
0	4	3515

We sample **1757** formulas from the above table (last row)

**500** heap graphs are generated per formula.

**878,500** formula/heap graph combinations.

Example input graph from dataset which is corrected predicted.  
(nested level 0, number of free variables 4)

Training / validation / testing sets: 6:2:2

Method	Top 1 Acc.	Top 10 Acc.
Concatenation features	0.05%	0.07%
Joint features	<b>91.5%</b>	<b>91.6%</b>

$\exists v_0. ls(arg0, v_0, \top)$   
 $* ls(arg1, arg3, \top)$   
 $* ls(arg3, arg3, \top)$   
 $* ls(v_0, v_0, \top)$

One can then use static program verification tools to determine whether the description is accurate and whether the program satisfies is memory safe.