

Memory abstraction 1

MPRI — Cours 2.6 “Interprétation abstraite :
application à la vérification et à l’analyse statique”

Xavier Rival

INRIA, ENS, CNRS

Dec, 10th. 2014

Overview of the lecture

So far, we have shown **numeric abstract domains**

- non relational: intervals, congruences...
- relational: polyhedra, octagons, ellipsoids...

- **How to deal with non purely numeric states ?**
- **How to reason about complex data-structures ?**

⇒ a **very broad topic**, and two lectures:

This lecture:

- **overview most common problems**
- discuss **arrays**, **strings**
- introduction to **shape analysis**

Next lecture: deeper study of a **family of shape analyses**

Assumptions

Programs can be viewed as **transition systems**:

- set of **control states**: \mathbb{L} (program points)
- set of **variables**: \mathbb{X} (all assumed globals)
- set of **values**: \mathbb{V} (**for now: \mathbb{V} consists of integers (or floats) only**)
- set of **memory states**: \mathbb{M} (**for now: $\mathbb{M} = \mathbb{X} \rightarrow \mathbb{V}$**)
- **error state**: Ω
- **states**: \mathbb{S}

$$\begin{aligned}\mathbb{S} &= \mathbb{L} \times \mathbb{M} \\ \mathbb{S}_\Omega &= \mathbb{S} \uplus \{\Omega\}\end{aligned}$$

- a program is described by a **transition relation**:

$$\rightarrow \subseteq \mathbb{S} \times \mathbb{S}_\Omega$$

Abstraction: described by a domain \mathbb{D}^\sharp and a concretization:

$$\gamma : (\mathbb{D}^\sharp, \sqsubseteq^\sharp) \longrightarrow (\mathcal{P}(\mathbb{S}), \subseteq)$$

Programs: syntax

We start with a **minimal language**, to be extended with arrays, strings, pointers...

A minimal imperative language

l	::=	l-values	
		x	($x \in \mathbb{X}$)
e	::=	expressions	
		c	($c \in \mathbb{V}$)
		l	(<i>lvalue</i>)
		$e \oplus e$	(<i>arithoperation, comparison</i>)
s	::=	statements	
		$l = e$	(assignment)
		$s; \dots s;$	(sequence)
		$\text{if}(e)\{s\}$	(condition)
		$\text{while}(e)\{s\}$	(loop)

Programs: semantics

We assume **classical definitions** for:

- **l-values**: $\llbracket l \rrbracket : \mathbb{M} \rightarrow \mathbb{X}$
- **expressions**: $\llbracket e \rrbracket : \mathbb{M} \rightarrow \mathbb{V}$
- **programs and statements**:
 - ▶ we assume a label **before each statement**
 - ▶ each statement defines a **set of transition** (\rightarrow)

We rely on the usual:

Reachable states semantics

The reachable states are computed as $\llbracket S \rrbracket_{\mathcal{R}} = \mathbf{lfp} F$ where

$$\begin{array}{lcl}
 F : \mathcal{P}(\mathbb{S}) & \longrightarrow & \mathcal{P}(\mathbb{S}) \\
 X & \longmapsto & \mathbb{S}_{\mathcal{I}} \cup \{s \in \mathbb{S} \mid \exists s' \in X, s' \rightarrow s\}
 \end{array}$$

Programs: semantics abstraction

We assume a **memory abstraction**:

- memory abstract domain $\mathbb{D}_{\text{mem}}^\#$
- concretization function $\gamma_{\text{mem}} : \mathbb{D}_{\text{mem}}^\# \rightarrow \mathcal{P}(\mathbb{M})$

Reachable states abstraction

We construct $\mathbb{D}^\# = \mathbb{L} \rightarrow \mathbb{D}_{\text{mem}}^\#$ and:

$$\begin{aligned} \gamma : \mathbb{D}^\# &\longrightarrow \mathcal{P}(\mathbb{S}) \\ X^\# &\longmapsto \{(\ell, m) \in \mathbb{S} \mid m \in \gamma_{\text{mem}}(X^\#(\ell))\} \end{aligned}$$

The whole question is how do we choose $\mathbb{D}_{\text{mem}}^\#$, γ_{mem} ...

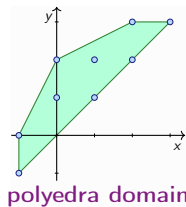
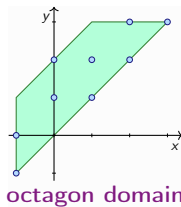
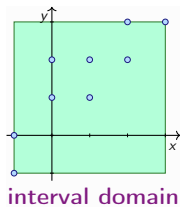
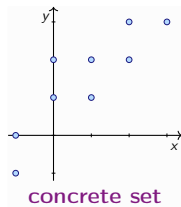
- previous lectures: \mathbb{X} is fixed and finite and, usually, \mathbb{V} is integers
- thus, $\mathbb{M} \equiv \mathbb{V}^n$

Abstraction of purely numeric memory states

Purely numeric case

- \mathbb{V} is a set of values of the same kind
- e.g., integers (\mathbb{Z}), machine integers ($\mathbb{Z} \cap [-2^{63}, 2^{63} - 1]$)...
- If the set of variables is fixed, we can use **any abstraction for \mathbb{V}^N**

Example: $N = 2$, $\mathbb{X} = \{x, y\}$



Heterogeneous memory states

In real life languages, there are many kinds of values:

- **scalars** (integers of various sizes, boolean, floating-point values)...
- **pointers, arrays**...

Heterogeneous memory states

- **types:** t_0, t_1, \dots
- **values:** $\mathbb{V} = \mathbb{V}_{t_0} \uplus \mathbb{V}_{t_1} \uplus \dots$
- finitely many **variables**; each has a **fixed type**: $\mathbb{X} = \mathbb{X}_{t_0} \uplus \mathbb{X}_{t_1} \uplus \dots$
- **memory states:**

$$\mathbb{M} = \mathbb{X}_{t_0} \rightarrow \mathbb{V}_{t_0} \times \mathbb{X}_{t_1} \rightarrow \mathbb{V}_{t_1} \dots$$

- At a later point, we will add **pointers**:
 t_0 **denotes pointers**, $\mathbb{V} = \dots \uplus \mathbb{V}_{\text{addr}}$
- For a moment, we let t_0 be integers, and t_1 be booleans

Heterogeneous memory states: non relational abstraction

Principle: compose abstractions for sets of memory states of each type

Non relational abstraction of heterogeneous memory states

- $\mathbb{M} \equiv \mathbb{M}_{t_0} \times \mathbb{M}_{t_1} \times \dots$ where $\mathbb{M}_{t_i} = \mathbb{X}_{t_i} \rightarrow \mathbb{V}_{t_i}$
- **Concretization function** (case with two types)

$$\begin{aligned} \gamma_{nr} : \mathcal{P}(\mathbb{M}_{t_0}) \times \mathcal{P}(\mathbb{M}_{t_1}) &\longrightarrow \mathcal{P}(\mathbb{M}) \\ (m_0^\sharp, m_1^\sharp) &\longmapsto \{(m_{t_0}, m_{t_1}) \mid \forall i, m_{t_i} \in \gamma_i(m_i^\sharp)\} \end{aligned}$$

Example: $\mathbb{V} = \mathbb{V}_{\text{int}} \uplus \mathbb{V}_{\text{bool}}$, thus, $\mathbb{M} = \mathbb{M}_{\text{int}} \times \mathbb{M}_{\text{bool}}$

Abstraction of $\mathcal{P}(\mathbb{X}_{\text{int}} \rightarrow \mathbb{V}_{\text{int}})$:

- intervals
- polyhedra...

Abstraction of $\mathcal{P}(\mathbb{X}_{\text{bool}} \rightarrow \mathbb{V}_{\text{bool}})$:

- lattice of boolean constants
- relational abstraction with BDDs

How about a relational analysis ?

Memory structures

- The definition $\mathbb{M} = \mathbb{X} \rightarrow \mathbb{V}$ is **too restrictive**
- It ignores many ways of organizing data in the memory states

Common structures (non exhaustive list)

- **Structures, records, tuples:**
sequences of cells accessed with fields
- **Arrays:** similar to structures; indexes are integers in $[0, n - 1]$
- **Pointers:**
numeric values corresponding to the address of a memory cell
- **Strings and buffers:**
blocks with a sequence of elements and a terminating element (e.g., *null character*)
- **Closures** (functional languages):
pointer to function code and (partial) list of arguments)

Specific properties to verify

Memory safety

Absence of memory errors (crashes, or undefined behaviors)

Pointer errors:

- Dereference of a **null pointer**
- Dereference of an **invalid pointer**

Access errors:

- Access to an array **out of its bounds**
- **Buffer overrun** (very commonly used for attacks)

Invariance properties

Data should not become corrupted (values or structures...)

Properties to verify: examples

A program closing a list of file descriptors

```
//l points to a list
c = l;
while(c ≠ NULL){
  close(c → FD);
  c = c → next;
}
```

Correctness properties

- memory safety
- l is supposed to store all file descriptors at all times
Will its structure be preserved ?
Yes, no breakage of a next link
- closure of all the descriptors

Examples of structure preservation properties

- Algorithms manipulating **trees**, **lists**...
- Libraries of algorithms on **balanced trees**
- **Not guaranteed by the language !**
e.g., balancing of Maps was wrong in the OCaml standard library...

Issues to consider in this lecture

- Propose a **concrete model**: expressive, intuitive...
- Abstract the **layout of memory states**
i.e., what is the structure of the data
- Abstract the **contents of data structures**
- Express **relations** among various elements
e.g., structural properties and properties of the contents of the structures
- Design **abstract interpretation algorithms**
 - ▶ transfer functions
 - ▶ widening

Outline

- 1 Towards memory properties
- 2 **Memory models**
 - Formalizing concrete memory states
 - Treatment of errors
 - Language semantics
- 3 Abstraction of arrays
- 4 Abstraction of strings and buffers
- 5 Basic pointer analyses
- 6 Three valued logic heap abstraction
- 7 Conclusion

A more realistic model

Not all memory cell corresponds to a variable

- a variable may correspond to **several cells**
- **heap allocated cells** correspond to no variable at all...

Environment + Heap

- **Addresses** are values: $\mathbb{V}_{\text{addr}} \subseteq \mathbb{V}$
- **Environments** $e \in \mathbb{E}$ map variables into their addresses
- **Heaps** ($h \in \mathbb{H}$) map addresses into values

$$\mathbb{E} = \mathbb{X} \rightarrow \mathbb{V}_{\text{addr}}$$

$$\mathbb{H} = \mathbb{V}_{\text{addr}} \rightarrow \mathbb{V}$$

h is actually only a partial function

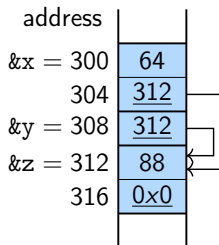
- **Memory states:** $\mathbb{M} = \mathbb{E} \times \mathbb{H}$

Example of a concrete memory state (variables)

- x and z are two list elements containing values 64 and 88, and where the former points to the latter
- y stores a pointer to z

Memory layout

(pointer values underlined)



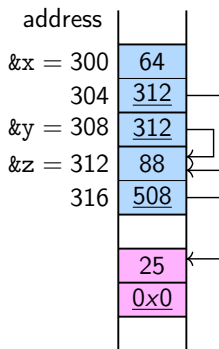
$$e : \begin{array}{l} x \mapsto 300 \\ y \mapsto 308 \\ z \mapsto 312 \end{array}$$

$$m : \begin{array}{l} 300 \mapsto 64 \\ 304 \mapsto 312 \\ 308 \mapsto 312 \\ 312 \mapsto 88 \\ 316 \mapsto 0 \end{array}$$

Example of a concrete memory state (variables + heap)

- same configuration
- + z points to a heap allocated list element (in purple)

Memory layout



$e :$	x	\mapsto	300
	y	\mapsto	308
	z	\mapsto	312
$m :$	300	\mapsto	64
	304	\mapsto	312
	308	\mapsto	312
	312	\mapsto	88
	316	\mapsto	508
	508	\mapsto	25
	512	\mapsto	0

Extending the language syntax

We start from the same language syntax and extend l-values:

l	::=	l-values	
		x	($x \in \mathbb{X}$)
		...	other kinds of l-values pointers, array dereference...
e	::=	expressions	
		c	($c \in \mathbb{V}$)
		l	(<i>lvalue</i>)
		e \oplus e	(<i>arithoperation, comparison</i>)
s	::=	statements	
		l = e	(assignment)
		s; ... s;	(sequence)
		if(e){s}	(condition)
		while(e){s}	(loop)

Extending the language semantics

Some slight modifications to the semantics of the initial language:

- **Values are addresses:** $\mathbb{V}_{\text{addr}} \subseteq \mathbb{V}$
- **L-values evaluate into addresses:** $\llbracket l \rrbracket : \mathbb{M} \rightarrow \mathbb{V}_{\text{addr}}$

$$\llbracket x \rrbracket(e, h) = e(x)$$

- **Semantics of expressions** $\llbracket e \rrbracket : \mathbb{M} \rightarrow \mathbb{V}_{\text{addr}}$, mostly unchanged

$$\llbracket l \rrbracket(e, h) = m(\llbracket l \rrbracket(e, h))$$

- **Semantics of assignment** $l_0 : l := e; l_1 : \dots$:

$$(l_0, e, h_0) \longrightarrow (l_1, e, h_1)$$

where

$$h_1 = h_0[\llbracket l \rrbracket(e, h_0) \leftarrow \llbracket e \rrbracket(e, h_0)]$$

Extensions of the symbolic model

Our model is still not quite realistic

- Memory cells do not all have the same **size**
- **Memory management algorithms** usually do not treat cells one by one, e.g., **malloc** returns a pointer to a *block* applying **free** to that pointer will dispose the *whole block*

Other refined models

- **Division** of the memory in **blocks** with a **base address** and a **size**
- **Division** of blocks into **cells** with a **size**
- Description of **fields** with an **offset**
- Description of **pointer values** with a **base address** and an **offset...**

For a **very formal** description of concrete memory states:
see **CompCert** project source files (Coq formalization)

Language semantics: program crash

- In an abnormal situation, **the program will crash**
- Advantage: very clear semantics
- Disadvantage (for the compiler designer): dynamic checks are required

Error state

- Ω denotes an **error configuration**
- Ω is a **blocking**: $\rightarrow \subseteq \mathbb{S} \times (\{\Omega\} \uplus \mathbb{S})$

OCaml:

- out-of-bound array access: Exception: `Invalid_argument`
"index out of bounds".
- no notion of a null pointer

Java:

- out-of-bound array access: exception
`java.lang.ArrayIndexOutOfBoundsException`

Language semantics: undefined behaviors

- The behavior of the program is **not specified** when an abnormal situation is encountered
- Advantage: easy implementation (often architecture driven)
- Disadvantage: unintuitive semantics, errors hard to reproduce

Modeling of undefined behavior

- Very hard to capture what a program operation may modify
- Abnormal situation at (l_0, m_0) m_0 such that
 $\forall m_1 \in \mathbb{M}, (l_0, m_0) \rightarrow (l_1, m_1)$
- **In C:**
Array out-of-bound accesses and dangling pointer dereferences
whereas a null-pointer dereference always result into a crash

Composite objects

How are contiguous blocks of information organized ?

Java objects, OCaml struct types

- sets of fields
- each field has its type
- **no assumption** on physical storage, **no pointer arithmetics**

C composite structures and unions

- **physical mapping** defined by the norm
- each field has a specified **size** and a specified **alignment**
- **union types / casts**:
implementations may allow several views

Pointers and records / structures / objects

- Our purpose is not to select a language for programming
- It is to remark salient language features, and their impact on abstractions

What kind of objects can be referred to by a pointer ?

Pointers only to records / structures / objects

- **Java**: only pointers to objects
- **OCaml**: only pointers to records, structures...

Pointers to fields

- **C**: pointers to any valid cell...

```
struct {int a; int b} x;  
int * y = &(x · b);
```


Pointer arithmetics

What kind of operations can be performed on a pointer ?

Classical pointer operations

- Pointer **dereference**:
 $\star p$ returns the contents of the cell pointed to by p
- **“Address of”** operator: $\&x$ returns the address of variable x
- Can be analyzed with **a rather coarse pointer model**
e.g., symbolic base + symbolic field

Arithmetics on pointers, requiring a more precise model

- **Addition of a numeric constant**:
 $p + n$: address contained in $p + n$ times the size of the type of p
Interaction with pointer casts...
- **Pointer subtraction**: returns a numeric offset

String operations

- Many **data-structures** can be handled in very different ways depending on the languages
- **Strings** are just one example

OCaml strings

- **Abstract type**: representation not part of the language definition
- **Type safe** implementation
 - ▶ no buffer overrun
 - ▶ exception for out of bound accesses
i.e., like arrays
- Most operations **generate new string structures**

C strings

- A **string** is an **array of characters (`char*`)** with a **terminal zero character**
- **Direct access** to string elements (array dereference)
- String copy operation **`strcpy(s, "foo_bar")`**:
 - ▶ copies "foo_bar" into s
 - ▶ **undefined behavior** if length of s < 7

Manual memory management

Allocation of unbounded memory space

- How are new memory blocks made available to the program ?
- How do old memory blocks get freed ?

OCaml memory management

- **Implicit allocation**
when declaring a new object
- **Garbage collection**: purely automatic process, that frees unreachable blocks

C memory management

- **Manual allocation**: **malloc** operation returns a pointer to a new block
- **Manual de-allocation**: **free** operation (block base address)

Manual memory management is not safe:

- **Memory leaks**: growing unreachable memory region; memory exhaustion
- **Dangling pointers** if freeing a block that is still referred to

Summary on the memory model

List of choices:

- **Clear error cases** or **undefined behaviors**
for analysis, a semantics with clear error cases is preferable
- **Composite objects**: structure fully exposed or not
- **Pointers to object fields**: allowed or not
- **Pointer arithmetic**: allowed or not
i.e., are pointer values symbolic values or numeric values
- **Memory management**: automatic or manual

We will generally assume a simple model,
unless considering specific features

Outline

- 1 Towards memory properties
- 2 Memory models
- 3 Abstraction of arrays**
 - A micro language for manipulating arrays
 - Verifying safety of array operations
 - Abstraction of array contents
 - Abstraction of array properties
- 4 Abstraction of strings and buffers
- 5 Basic pointer analyses
- 6 Three valued logic heap abstraction

Programs: extension with arrays

Extension of the syntax:

l	$::=$	l-valules	
		\dots	previous constructions
		$x[e]$	cell of array x
\dots	$::=$	\dots	the rest is unchanged

Extension of the semantics:

- if x is an array variable, and corresponds to an array of length N , we have N cells corresponding to it, with addresses

$$\{e(x) + 0, e(x) + s, \dots, e(x) + (N - 1)s\}$$

where s is the size of an array cell (e.g., 8 bytes for a 64-bit int)

- evaluation of an array cell read:**

$$\llbracket x[e] \rrbracket(e, h) = \begin{cases} e(x) + is & \text{if } \llbracket e \rrbracket(e, h) = i \in [0, N - 1] \\ \Omega & \text{otherwise} \end{cases}$$

Example

```
// a is an integer array of length n
bool s;
do{
  s = false;
  for(int i = 0; i < n - 1; i = i + 1){
    if(a[i] < a[i + 1]){
      swap(a[i], a[i + 1]);
      s = true;
    }
  }
} while(s);
```

Properties to verify by static analysis

- 1 **Safety property:** the program will not crash (no index out of bound)
- 2 **Contents property:** if the values in the array are in $[0, 100]$ before, they are also in that range after
- 3 **Global array property:** at the end, the array is sorted

Outline

- 1 Towards memory properties
- 2 Memory models
- 3 Abstraction of arrays**
 - A micro language for manipulating arrays
 - Verifying safety of array operations
 - Abstraction of array contents
 - Abstraction of array properties
- 4 Abstraction of strings and buffers
- 5 Basic pointer analyses
- 6 Three valued logic heap abstraction

Expressing correctness of array operations

Goal of the analysis: establish safety

Prove the absence of runtime error due to array reads / writes
i.e., **that no Ω will ever arise**

Safety verification:

- At label ℓ , the analysis computes a **local abstraction of the set of reachable memory states** $\Phi^\sharp(\ell)$
- If a statement at label ℓ performs array read or write operation $\mathbf{x}[e]$, where \mathbf{x} is an array of length n , the analysis simply needs to establish
$$\forall m \in \gamma_{\text{mem}}(\Phi^\sharp(\ell)), \llbracket e \rrbracket(m) \in [0, n - 1]$$
- In many cases, this can be done with an **interval abstraction**
... but not always (Exercise: when would it not be enough ?)

For now, we ignore the contents of the array
(Exercise: when does this fail ?)

Verifying correctness of array operations

Case where intervals are enough:

```
// x array of length 40
int i = 0;
while(i < 40){
    printf("%d;", x[i]);
    i = i + 1;
}
```

- interval analysis establishes that $i \in [0; 39]$ at the loop head
- this allows the verification of the code

Case where intervals cannot represent precise enough invariants:

```
// x array of length 40
int i, j;
if(0 ≤ i && i < j)
    if(j < 41)
        printf("%d;", x[i]);
```

- in the concrete, $i \in [0, 39]$ at the array access point
- to establish this in the abstract, after the first test, relation $i < j$ need be represented
- e.g., octagon abstract domain

Outline

- 1 Towards memory properties
- 2 Memory models
- 3 Abstraction of arrays**
 - A micro language for manipulating arrays
 - Verifying safety of array operations
 - Abstraction of array contents
 - Abstraction of array properties
- 4 Abstraction of strings and buffers
- 5 Basic pointer analyses
- 6 Three valued logic heap abstraction

Elementwise abstraction

Goal of the analysis: abstract contents

Inferring invariants about the **contents** of the array

- e.g., that the values in the array **are in a given range**
- e.g., in order to verify the **safety of $x[y[i + j] + k]$** or $y = n/x[i]$

Assumption:

- **One array t** , of **known, fixed length n** (element size s)
- Scalar variables x_0, x_1, \dots, x_{m-1}

Concrete memory cell addresses:

$$\mathbb{V}_{\text{addr}} = \{\&x_0, \dots, \&x_{m-1}\} \cup \{\&\bar{t}, \&\bar{t} + 1 \cdot s, \dots, \&\bar{t} + (n - 1) \cdot s\}$$

Elementwise abstraction

- **Each** concrete cell is **mapped into one abstract cell**
- $\mathbb{D}^\#$ should simply be an **abstraction of $\mathcal{P}(\mathbb{V}^{m+n})$**

Array smashing abstraction: abstraction into one cell

The elementwise abstraction is **too costly**:

- **high number of abstract cells** if the arrays are big
- **will not work** if the size of arrays is **not known statically**

Alternative: **use fewer abstract cells**, e.g., **a single cell**

Assumption: m scalar variables, \bar{t} array of length n

Array smashing

- All cells of the array are mapped into **one abstract cell** \bar{t}
- **Abstract cells**: $\mathbb{C}^\# = \{\&x_0, \dots, \&x_{m-1}\} \cup \{\&\bar{t}\}$
- $\mathbb{D}^\#$ should simply be an **abstraction of** $\mathcal{P}(\mathbb{V}^{m+1})$

This also works **if the size of the array is not known statically**:

```
int n = ...;
int t[n];
```

The contents of t is represented using one abstract cell whatever the value of n

Array smashing abstraction

Definition

- **Abstract domain** $\mathcal{P}(\mathbb{C}^\# \rightarrow \mathcal{P}(\mathbb{V}))$
- **Abstraction function:**

$$\alpha_{\text{smash}}(H) = \left\{ \begin{array}{l} \&x_i \mapsto \{h(x_i)\} \\ \&\bar{t} \mapsto \{h(\&t + 0), \dots, h(\&t + n - 1)\} \end{array} \mid h \in H \right\}$$

Example:

- No variable, array of length 2
- **Set of concrete states:**

$$\left\{ \begin{array}{l} t[0] \mapsto 0 \\ t[1] \mapsto 10 \end{array} \right\}, \quad \left\{ \begin{array}{l} t[0] \mapsto 2 \\ t[1] \mapsto 11 \end{array} \right\}, \quad \left\{ \begin{array}{l} t[0] \mapsto 1 \\ t[1] \mapsto 12 \end{array} \right\}$$

- **Abstract state**, using interval abstraction: $\&\bar{t} \mapsto [0, 12]$

Weak updates: an imprecision in the analysis

Assumptions:

- **Smashing abstraction**, with the **interval abstract domain**
- Array t is supposed **of known length** $n \geq 2$
- We consider statement $\ell_0 : t[i] = 0; \ell_1$
- Given m_0^\sharp , using intervals to describe a set of states at ℓ , we wish to compute an over-approximation m_1^\sharp of

$$\{m_1 \mid \exists m_0 \in \gamma_{\text{mem}}(m_0^\sharp), (\ell_0, m_0) \rightarrow (\ell_1, m_1)\}$$
- **Abstract pre-condition:** $m_0^\sharp(\&i) = [0, 0], m_0^\sharp(\&\bar{t}) = [a, b]$

Post-condition:

- in the **concrete** level:

$$\begin{cases} \&t + 0 &\longmapsto & 0 & \text{(cell just modified)} \\ \&t + 1 &\longmapsto & v & \text{where } v \in [a, b] \text{ (cell not modified)} \end{cases}$$
- in the **abstract** level, we **only lose precision**:

$$\&\bar{t} \longmapsto [0, 0] \sqcup [a, b] = [\min(a, 0), \max(b, 0)]$$

Weak updates

Summary:

- i was known very precisely
- $\&\bar{t}$ **stands for several concrete cells**
- The assignment will modify only **one cell**
the others will keep their old value
- The abstraction cannot distinguish unmodified values from the modified cell
- As a consequence, the range for $\&\bar{t}$ **may only grow**

Weak updates

- It would only be worse if the value of i was not known precisely
- This is a significant loss in precision
- This is a **limitation of all smashing analyses**

Weak updates and strong updates

Definitions

- **Strong update:** modified abstract cell **fully materialized**, and **old value fully discarded**
- **Weak update:** modified abstract cell **not fully materialized**, and **new value “joined” with old values**

In the case of $t[i] := e$, weak updates may arise in the following cases:

- using a **smashing abstraction**:
 \bar{t} denotes several concrete cells; only one gets modified, so we must keep old values
- using a **pointwise abstraction**, if $m_0^\sharp(i) = [i, i']$ where $i < i'$:
 - ▶ one cell in $\{\&t + i \cdot s, \dots, \&t + i' \cdot s\}$ gets modified
 - ▶ the other cells in that set remain the same
 - ▶ so we must also keep old values

Weak updates and strong updates: example

```
//x uninitialized array of length n
int i = 0;
while(i < n){
    x[i] = 0;
    i = i + 1;
}
```

Pointwise abstraction:

- initially $\forall i, m^\sharp(\&t + i \cdot s) = \top$
- if loop unrolled completely, at the end, $\forall i, m^\sharp(\&t + i \cdot s) = [0, 0]$
- weak updates, if the loop is not unrolled; then, at the end $\forall i, m^\sharp(\&t + i \cdot s) = \top$

Smashing abstraction:

- initially $m^\sharp(\bar{t}) = \top$
- weak updates at each step (whatever the unrolling that is performed); at the end: $m^\sharp(\bar{t}) = \top$

- Weak updates may cause a **serious loss of precision**
- Workaround ahead: **more complex array abstractions** may help

Other forms of array smashing

- Smashing does not have to affect the whole array
- Efficient smashing strategies can be found

Segment smashing:

- abstraction of the array cells into $\{\bar{t}_0, \dots, \bar{t}_{k-1}\}$ where \bar{t}_i corresponds to **a segment of the array**
- useful when sub-segments have interesting properties
- **issue**: determine the segment by analysis

Modulo smashing:

- abstraction of the array cells into $\{\bar{t}_0, \dots, \bar{t}_{k-1}\}$ where \bar{t}_i corresponds to **a repeating set of offsets** $\{\&\bar{t} + k \cdot i \cdot s \mid k \cdot i < n\}$
- useful for arrays of structures
- **issue**: determine k by analysis

Outline

- 1 Towards memory properties
- 2 Memory models
- 3 Abstraction of arrays**
 - A micro language for manipulating arrays
 - Verifying safety of array operations
 - Abstraction of array contents
 - Abstraction of array properties
- 4 Abstraction of strings and buffers
- 5 Basic pointer analyses
- 6 Three valued logic heap abstraction

Example array properties

Goal of the analysis: precisely abstract contents

Discover non trivial properties of **array regions**

- Initialization to a constant (e.g., 0)
- Sortedness

An array initialization loop:

```
// t integer array of length n
int i = 0;
while(i < n){
    t[i] = 0;
    i = i + 1;
}
```

Sketch of a hand proof:

- At iteration i , $i = i$ and the segment $t[0], \dots, t[i - 1]$ is initialized
- At the loop exit, $i = n$ and the whole array is initialized

We need to express properties on segments;
otherwise the proof cannot be completed

Array segment properties

An array initialization loop:

```
// t integer array of length n
int i = 0;
while(i < n){
    t[i] = 0;
    i = i + 1;
}
```

Concrete state after 6 iterations:

$i = 6$

t	0	0	0	0	0	0	?	?	?	?
---	---	---	---	---	---	---	---	---	---	---

Corresponding abstract state:

$i \in [1, 10]$

t	$\text{zero}_t(0, i - 1)$	\top
---	---------------------------	--------

Array segment predicates

Definition

An **array segment predicate** is an abstract predicate that describes the contents of a contiguous series of cells in the array, such as:

- **Initialization**: $\text{zero}_t(i, j)$ iff t initialized to 0 between i and j
- **Sortedness**: $\text{sort}_t(i, j)$ iff t sorted between i and j

Examples:

- array satisfying $\text{zero}_t(2, 6)$:

$$i = 6$$

t	8	2	0	0	0	0	0	0	10	3
---	---	---	---	---	---	---	---	---	----	---

- array satisfying $\text{sort}_t(1, 4)$ and $\text{sort}_t(6, 8)$:

$$i = 6$$

t	8	2	5	6	8	11	1	2	3	2
---	---	---	---	---	---	----	---	---	---	---

Composing sortedness predicates

As part of the proof, predicates need be composed

$$\begin{aligned}
 \text{zero}_t(i, j) \wedge \text{zero}_{\bar{t}}(j + 1, k) &\Rightarrow \text{zero}_t(i, k) \\
 \text{zero}_t(i, j) \wedge t[j + 1] = 0 &\Rightarrow \text{zero}_t(i, j + 1) \\
 \text{sort}_t(i, j) \wedge \text{sort}_{\bar{t}}(j + 1, k) &\not\Rightarrow \text{sort}_t(i, k) \\
 t[j] \leq t[j + 1] \wedge \text{sort}_t(i, j) \wedge \text{sort}_{\bar{t}}(j + 1, k) &\Rightarrow \text{sort}_t(i, k)
 \end{aligned}$$

- **counter example** for the third line: for $[0; 3; 9; 2; 4; 8]$, we have:

$$\text{sort}_t(0, 2) \wedge \text{sort}_t(3, 5) \quad \text{but not} \quad \text{sort}_t(0, 5)$$

Another sortedness predicate: $\text{sort}_t(i, j, \text{min}, \text{max})$

$$B \leq C \wedge \text{sort}_t(i, j, A, B) \wedge \text{sort}_{\bar{t}}(j + 1, k, C, D) \Rightarrow \text{sort}_t(i, k, A, D)$$

Analysis operators (for predicate **zero**)

Assignment transfer function:

- 1 Identify segments that may be modified
- 2 If a single segment is impacted, split it
- 3 Do a strong update


$$\begin{array}{l} \mathbf{zero}_t(0, n) \wedge 0 \leq i < n \xrightarrow{t[i]=?} \mathbf{zero}_t(0, i-1) \wedge \mathbf{zero}_t(i+1, n) \wedge 0 \leq i < n \\ \top \wedge 0 \leq i < n \xrightarrow{t[i]=0} \mathbf{zero}_t(i, i) \wedge 0 \leq i < n \end{array}$$

Abstract join operator: generalizes bounds

$$(\top \wedge i = 0 < n) \sqcup^{\sharp} (\mathbf{zero}_t(0, 0) \wedge i = 1 < n) = (\mathbf{zero}_t(0, i-1) \wedge 0 \leq i < n)$$

Array analysis: example

```
// t integer array of length  $n > 0$ 
```

```
  t  i T
```

The diagram shows a horizontal rectangle representing an array. Inside the rectangle, the letter 'T' is centered. To the left of the rectangle is the variable 't', and to the right is the variable 'i' followed by a space and the letter 'T'.

```
int i = 0;
```

```
  t  i T
```


The diagram shows a horizontal rectangle representing an array. Inside the rectangle, the letter 'T' is centered. To the left of the rectangle is the variable 't', and to the right is the variable 'i' followed by a space and the letter 'T'.

```
while(i < n){
```

```
  t  i T
```


The diagram shows a horizontal rectangle representing an array. Inside the rectangle, the letter 'T' is centered. To the left of the rectangle is the variable 't', and to the right is the variable 'i' followed by a space and the letter 'T'.

```
  t[i] = 0;
```

```
  t  i T
```

The diagram shows a horizontal rectangle representing an array. Inside the rectangle, the letter 'T' is centered. To the left of the rectangle is the variable 't', and to the right is the variable 'i' followed by a space and the letter 'T'.

```
  i = i + 1;
```

```
  t  i T
```

The diagram shows a horizontal rectangle representing an array. Inside the rectangle, the letter 'T' is centered. To the left of the rectangle is the variable 't', and to the right is the variable 'i' followed by a space and the letter 'T'.

```
}
```

```
  t  i T
```

The diagram shows a horizontal rectangle representing an array. Inside the rectangle, the letter 'T' is centered. To the left of the rectangle is the variable 't', and to the right is the variable 'i' followed by a space and the letter 'T'.

Array analysis: example

```
// t integer array of length  $n > 0$ 
```

```
t  i T
```

The diagram shows a horizontal rectangle representing an array. Inside the rectangle, there is a large capital letter 'T' centered. To the left of the rectangle is the variable 't', and to the right is the variable 'i' followed by the value 'T'.

```
int i = 0;
```

```
t  i [0,0]
```

The diagram shows a horizontal rectangle representing an array. Inside the rectangle, there is a large capital letter 'T' centered. To the left of the rectangle is the variable 't', and to the right is the variable 'i' followed by the value '[0,0]'.

```
while(i < n){
```

```
t  i T
```

The diagram shows a horizontal rectangle representing an array. Inside the rectangle, there is a large capital letter 'T' centered. To the left of the rectangle is the variable 't', and to the right is the variable 'i' followed by the value 'T'.

```
t[i] = 0;
```

```
t  i T
```

The diagram shows a horizontal rectangle representing an array. Inside the rectangle, there is a large capital letter 'T' centered. To the left of the rectangle is the variable 't', and to the right is the variable 'i' followed by the value 'T'.

```
i = i + 1;
```

```
t  i T
```

The diagram shows a horizontal rectangle representing an array. Inside the rectangle, there is a large capital letter 'T' centered. To the left of the rectangle is the variable 't', and to the right is the variable 'i' followed by the value 'T'.

```
}
```

```
t  i T
```

The diagram shows a horizontal rectangle representing an array. Inside the rectangle, there is a large capital letter 'T' centered. To the left of the rectangle is the variable 't', and to the right is the variable 'i' followed by the value 'T'.

Array analysis: example

```
// t integer array of length  $n > 0$ 
```

t		i	T
---	---	---	---

```
int i = 0;
```

t		i	[0, 0]
---	---	---	--------

```
while(i < n){
```

t		i	[0, 0]
---	---	---	--------

```
  t[i] = 0;
```

t		i	T
---	---	---	---

```
  i = i + 1;
```

t		i	T
---	---	---	---

```
}
```

t		i	T
---	---	---	---

Array analysis: example

```
// t integer array of length  $n > 0$ 
```

t	T	i	T
---	---	---	---

```
int i = 0;
```

t	T	i	[0, 0]
---	---	---	--------

```
while(i < n){
```

t	T	i	[0, 0]
---	---	---	--------

```
  t[i] = 0;
```

t	zero _t (0, 1)	T	i	[0, 0]
---	--------------------------	---	---	--------

```
  i = i + 1;
```

t	T	i	T
---	---	---	---

```
}
```

t	T	i	T
---	---	---	---

Array analysis: example

// t integer array of length $n > 0$

t  i T

The diagram shows a horizontal rectangle representing an array. The letter 'T' is centered inside the rectangle.

int i = 0;

t  i [0, 0]

The diagram shows a horizontal rectangle representing an array. The letter 'T' is centered inside the rectangle.

while(i < n){

t  i [0, 0]

The diagram shows a horizontal rectangle representing an array. The letter 'T' is centered inside the rectangle.

t[i] = 0;

t  i [0, 0]

The diagram shows a horizontal rectangle representing an array. The first part of the rectangle, from the left edge to approximately one-fifth of the way across, is highlighted in yellow and contains the text $\text{zero}_{\bar{i}}(0, 1)$. The rest of the rectangle is grey and contains the letter 'T'.

i = i + 1;

t  i [1, 1]

The diagram shows a horizontal rectangle representing an array. The first part of the rectangle, from the left edge to approximately one-fifth of the way across, is highlighted in yellow and contains the text $\text{zero}_{\bar{i}}(0, 1)$. The rest of the rectangle is grey and contains the letter 'T'.

}

t  i T

The diagram shows a horizontal rectangle representing an array. The letter 'T' is centered inside the rectangle.

Array analysis: example

```
// t integer array of length  $n > 0$ 
```

t	T	i	T
---	---	---	---

```
int i = 0;
```

t	zero _t (0, i - 1)	T	i	[0, 1]
---	------------------------------	---	---	--------

```
while(i < n){
```

t	T	i	[0, 0]
---	---	---	--------

```
  t[i] = 0;
```

t	zero _t (0, 1)	T	i	[0, 0]
---	--------------------------	---	---	--------

```
  i = i + 1;
```

t	zero _t (0, 1)	T	i	[1, 1]
---	--------------------------	---	---	--------

```
}
```

t	T	i	T
---	---	---	---

Array analysis: example

// t integer array of length $n > 0$

t  i T

int i = 0;

t  i [0, 1]

while(i < n){

t  i [0, 1]

t[i] = 0;

t  i [0, 0]

i = i + 1;

t  i [1, 1]

}

t  i T

Array analysis: example

// t integer array of length $n > 0$

t  i \top

int i = 0;

t  i [0, 1]

while(i < n){

t  i [0, 1]

t[i] = 0;

t  i [0, 1]

i = i + 1;

t  i [1, 1]

}

t  i \top

Array analysis: example

```
// t integer array of length  $n > 0$ 
```

t	T	i	T
---	---	---	---

```
int i = 0;
```

t	$\text{zero}_{\bar{\epsilon}}(0, i - 1)$	T	i	$[0, 1]$
---	--	---	---	----------

```
while(i < n){
```

t	$\text{zero}_{\bar{\epsilon}}(0, i - 1)$	T	i	$[0, 1]$
---	--	---	---	----------

```
  t[i] = 0;
```

t	$\text{zero}_{\bar{\epsilon}}(0, i)$	T	i	$[0, 1]$
---	--------------------------------------	---	---	----------

```
  i = i + 1;
```

t	$\text{zero}_{\bar{\epsilon}}(0, i - 1)$	T	i	$[1, 2]$
---	--	---	---	----------

```
}
```

t	T	i	T
---	---	---	---

Array analysis: example


// t integer array of length $n > 0$

t  i T

int i = 0;

t  i [0, n]

while(i < n){

t  i [0, 1]

t[i] = 0;

t  i [0, 1]

i = i + 1;

t  i [1, 2]

}

t  i T

Array analysis: example

// t integer array of length $n > 0$

t  i T

int i = 0;

t  i [0, n]

while(i < n){

t  i [0, n - 1]

t[i] = 0;

t  i [0, 1]

i = i + 1;

t  i [1, 2]

}

t  i T

Array analysis: example

// t integer array of length $n > 0$

t  i T

int i = 0;

t  i [0, n]

while(i < n){

t  i [0, n - 1]

t[i] = 0;

t  i [0, n - 1]

i = i + 1;

t  i [1, 2]

}

t  i T

Array analysis: example

// t integer array of length $n > 0$

t  i T

int i = 0;

t  i [0, n]

while(i < n){

t  i [0, n - 1]

t[i] = 0;

t  i [0, n - 1]

i = i + 1;

t  i [1, n]

}

t  i T

Array analysis: example

// t integer array of length $n > 0$

t  i T

int i = 0;

t  i [0, n]

while(i < n){

t  i [0, n - 1]

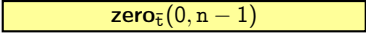
t[i] = 0;

t  i [0, n - 1]

i = i + 1;

t  i [1, n]

}

t  i [n, n]

Partitioning of arrays

Array partitions

A partition of an array t of length n is a sequence $\mathcal{P} = \{e_0, \dots, e_k\}$ of symbolic expressions where

- e_i denotes the lower (resp., upper) bound of element i (resp. $i - 1$) of the partition
- e_0 should be equal to 0 (and e_k to n)

Example:

- set of four **concrete states**:

$$\left\{ \begin{array}{ll} i = 1 & [0, 4, 1, 2, 3, 5] \\ i = 2 & [0, 1, 5, 2, 3, 4] \end{array} \right. \quad \begin{array}{ll} i = 3 & [2, 2, 4, 5, 1, 8] \\ i = 5 & [0, 2, 4, 6, 7, 9] \end{array}$$

- **partition**: $\{0, i + 1, 6\}$
- note that the array is always
 - ▶ sorted between 0 and i
 - ▶ sorted between $i + 1$ and 5

Abstraction based on array partitions

Segment and array abstraction

An array segmentation is given by a partition $\mathcal{P} = \{e_0, \dots, e_k\}$ and a set of abstract properties $\{P_0, \dots, P_{k-1}\}$.

Its concretization is the set of memory states $m = (e, \hat{h})$ such that

$\forall i, [t[v], t[v+1], \dots, t[w-1]]$ satisfies P_i , where $\begin{cases} v = \llbracket e_i \rrbracket(m) \\ w = \llbracket e_{i+1} \rrbracket(m) \end{cases}$

- **Partitions can be:**
 - ▶ **static**, i.e., pre-computed by another analysis [HP'08]
 - ▶ **dynamic**, i.e., computed as part of the analysis [CCL'11]
(more complex abstract domain structure with partitions *and* predicates)
- **Example:** array initialization

Outline

- 1 Towards memory properties
- 2 Memory models
- 3 Abstraction of arrays
- 4 Abstraction of strings and buffers**
 - A micro-language with strings
 - Abstraction
- 5 Basic pointer analyses
- 6 Three valued logic heap abstraction
- 7 Conclusion

Strings in programming languages

- In **high-level programming languages**:
 - ▶ **high-level** API, like OCaml `String` module or Java `String` classes
 - ▶ a set of **exceptions** in case of an invalid operation
 - ▶ **no security** risk in case of a crash
- In **C**:
 - ▶ **arrays of characters**
 - ▶ integration in other structures **with no protection**
 - ▶ **direct access**, with **no protection**

We focus on the case of languages with strings *à la C*

Programs: syntax and semantics

We extend our simple language with strings...

Encoding of strings in C

- **Strings** are represented by **character arrays**, with a **terminating 0**
- Only characters to the first zero are meaningful
- Example of a **string buffer** of length 10 containing string "hello"

'h'	'e'	'l'	'l'	'o'	'/'0'	'b'	'/'0'	'a'	'x'
-----	-----	-----	-----	-----	-------	-----	-------	-----	-----

Thus, **the language is essentially the same as for arrays**:

- data-types remain the same; we include a **char** type;
- expressions and l-values remain the same too
- we consider a set of **string operations** (typically, library functions)

Programs: string operations

String operations

- `strcpy(char * d, char * s)`: copies `s` into `d`, including terminating `0`, provided there is enough space (unspecified otherwise)
- `strncpy(char * d, char * s, int n)`: copies exactly `n` characters at most, from `s` into `d`
- `printf`: interprets `"%s"` as a string placeholder; displays up to the terminating `0` (unspecified if there is none)

```
char q[2];
char s[2];
char t[4];
strcpy(t, "bon");
strncpy(s, t, 2);
strcpy(q, s);
printf("nres: %s/n", q);
```

Result ?

- **not fully defined**
- depends on the **order** of memory blocks in memory...

Abstraction of string buffers

Goal of static analysis

Prove the absence of runtime errors in string buffer operations

Such errors could:

- cause **abrupt crashes** (segmentation fault) or undefined behaviors
- make **exploits** possible (e.g., by overwriting other program data)

We remark that:

- the **positions of “zero” characters** matters
- the **value of the other characters** usually does not matter
exception: cases where the program decides what to do depending on non zero characters, and where that impacts the error behavior of the program

Numeric abstraction of strings

String characters abstractions

We consider the character abstraction below:

$$\begin{aligned} \phi : \emptyset &\mapsto \emptyset & \phi : c &\mapsto '?' \\ \phi : c_0 \cdots c_{n-1} &\mapsto \phi(c_0) \cdots \phi(c_{n-1}) \\ \alpha_{\text{string}} : \mathcal{S} &\mapsto \{\phi(s) \mid s \in \mathcal{S}\} \end{aligned}$$

- α_{string} abstracts unneeded characters information

Numerical abstraction

We consider memory states that comprise only one string buffer t . We can abstract each such state using two numbers

- t_n : size of buffer t
- t_z : position of the first 0 in t if any (otherwise, we let $t_z = t_n$)

Abstraction of string buffers

We consider a program with integer variables $\mathbb{X}_{\text{int}} = \{x, y, \dots\}$ and string buffer variables $\mathbb{X}_{\text{buf}} = \{t, u, \dots\}$

Abstract domain

- We let $\mathbb{X}' = \mathbb{X}_{\text{int}} \uplus \{t_n, t_z, u_n, u_z, \dots\}$
- Each memory state m gets abstracted into a state $m' = \mathbf{abs}(m)$ over \mathbb{X}'
- Given an abstract domain $(\mathbb{D}_{\text{num}}^{\#}, \sqsubseteq_{\text{num}})$ of $\mathcal{P}(\mathbb{X}' \rightarrow \mathbb{Z})$, we can build an abstraction of $(\mathcal{P}(\mathbb{M}), \subseteq)$:

$$\begin{aligned} \gamma_{\text{buf}} : \quad \mathbb{D}_{\text{num}}^{\#} &\longrightarrow \mathcal{P}(\mathbb{M}) \\ X^{\#} &\longmapsto \{m \in \mathbb{M} \mid \mathbf{abs}(m) \in \gamma_{\text{num}}(X^{\#})\} \end{aligned}$$

Typical choice: polyhedra

Example

- **Example:** abstraction of

'h'	'e'	'l'	'l'	'o'	'/'	'o'	'b'	'/'	'a'	'x'
-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----

 into $t_n = 10, t_z = 5$
- **Practical implementation:**
 - ▶ either as a classical static analysis
 - ▶ or using a **transformation into an integer program**
- **Code transformation approach:**

<pre>char q[2]; char s[2]; char t[4]; strcpy(t, "bon"); strncpy(s, t, 2); strcpy(q, s); printf("nres: %s/n", q);</pre>	}	→	<pre>q_n = 2; s_n = 2; t_n = 2; t_z = 3; if(t_z < 2){s_z = t_z;} else if(s_z < t_n){s_z = s_n} assert(s_z < q_n); q_z = s_z; assert(q_z < q_n);</pre>
--	---	---	---

Outline

- 1 Towards memory properties
- 2 Memory models
- 3 Abstraction of arrays
- 4 Abstraction of strings and buffers
- 5 Basic pointer analyses**
 - A micro-language with pointers
- 6 Three valued logic heap abstraction
- 7 Conclusion

Programs with pointers: syntax

Syntax extension: quite a few additional constructions

l	::= l-valules	
	x	($x \in \mathbb{X}$)
	...	
	*e	pointer dereference
	l · f	field read
e	::= expressions	
	l	
	...	
	&l	"address of" operator
s	::= statements	
	...	
	x = malloc(c)	allocation of <i>c</i> bytes
	free(x)	deallocation of the block pointed to by <i>x</i>

We do not consider **pointer arithmetics here**

Programs with pointers: semantics

Case of l-values:

$$\llbracket \mathbf{x} \rrbracket(e, h) = e(\mathbf{x})$$

$$\llbracket *e \rrbracket(e, h) = \begin{cases} h(\llbracket e \rrbracket(e, h)) & \text{if } \llbracket e \rrbracket(e, h) \neq 0 \wedge \llbracket e \rrbracket(e, h) \in \mathbf{Dom}(h) \\ \Omega & \text{otherwise} \end{cases}$$

$$\llbracket \mathbf{l} \cdot \mathbf{f} \rrbracket(e, heap) = \llbracket \mathbf{l} \rrbracket(e, h) + \mathbf{offset}(\mathbf{f}) \text{ (numeric offset)}$$

Case of expressions:

$$\llbracket \mathbf{l} \rrbracket(e, heap) = h(\llbracket \mathbf{l} \rrbracket(e, heap))$$

$$\llbracket \&\mathbf{l} \rrbracket(e, heap) = \llbracket \mathbf{l} \rrbracket(e, heap)$$

Case of statements:

- **memory allocation** $\mathbf{x} = \mathbf{malloc}(c)$: $(e, h) \rightarrow (e, h')$ where $h' = h[e(\mathbf{x}) \leftarrow k] \uplus \{k \mapsto v_k, k+1 \mapsto v_{k+1}, \dots, k+c-1 \mapsto v_{k+c-1}\}$ and $k, \dots, k+c-1$ are fresh in h
- **memory deallocation** $\mathbf{free}(\mathbf{x})$: $(e, h) \rightarrow (e, h')$ where $k = e(\mathbf{x})$ and $h' = h \uplus \{k \mapsto v_k, k+1 \mapsto v_{k+1}, \dots, k+c-1 \mapsto v_{k+c-1}\}$

Pointer non relational abstraction: null pointers

The dereference of a null pointer will cause programs to crash

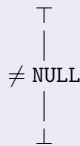
We go back to the **non relational abstraction of heterogeneous states**

- $\mathbb{V} = \mathbb{V}_{\text{addr}} \uplus \mathbb{V}_{\text{int}}$, $\mathbb{X} = \mathbb{X}_{\text{addr}} \uplus \mathbb{X}_{\text{int}}$
- we apply a non relational abstraction to pointer variables, based on $\mathbb{D}_{\text{addr}}^{\#}$ and $\gamma_{\text{addr}} : \mathbb{D}_{\text{addr}}^{\#} \rightarrow \mathcal{P}(\mathbb{V}_{\text{addr}})$

Null pointer analysis

Abstract lattice for addresses:

- $\gamma_{\text{addr}}(\perp) = \emptyset$
- $\gamma_{\text{addr}}(\top) = \mathbb{V}_{\text{addr}}$
- $\gamma_{\text{addr}}(\neq \text{NULL}) = \mathbb{V}_{\text{addr}} \setminus \{0\}$



- very lightweight, can typically resolve rather trivial cases
- useful for C, but also for Java

Pointer non relational abstraction: dangling pointers

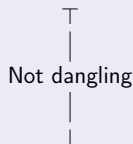
The dereference of a null pointer will cause programs to crash

This requires a **similar abstraction**:

Null pointer analysis

Abstract lattice for addresses:

- $\gamma_{\text{addr}}(\perp) = \emptyset$
- $\gamma_{\text{addr}}(\top) = \mathbb{V}_{\text{addr}} \times \mathbb{H}$
- $\gamma_{\text{addr}}(\text{Not dangling}) = \{(v, h) \mid h \in \mathbb{H} \wedge v \in \text{Dom}(h)\}$



- very lightweight, can typically resolve rather trivial cases
- useful for C
- in Java, superseded by the requirement that any variable be initialized

Pointer non relational abstraction: pointer aliasing

Determine where a pointer may store a reference to

Very useful to **support client analyses**:

```

1 : int x, y;
2 : int * p;
3 : y = 9;
4 : p = &x;
5 : *p = 0;

```

- what is the final value for x ?
0, since **it is modified at line 5**...
- what is the final value for x ?
0, since **it is not modified at line 5**...

Basic pointer abstraction

- We assume a set of **abstract memory locations** \mathbb{A}^\sharp is fixed:

$$\mathbb{A}^\sharp = \{\&x, \&y, \dots, \&t, a_0, a_1, \dots, a_N\}$$

- All concrete addresses are abstracted into \mathbb{A}^\sharp
- A pointer value is abstracted by the abstraction of the addresses it may point to (example, for p: $\{\&x\}$)

Pointer aliasing based on equivalence on access paths

Aliasing relation

Given $m = (e, h)$, pointers p and q are **aliases** iff $h(e(p)) = h(e(q))$

Abstraction to infer pointer aliasing properties

- An **access path** describes a sequence of operations to compute an l-value (i.e., an address); e.g.:

$$a ::= x \mid a \cdot f \mid *a$$

- An **abstraction for aliasing** is an over-approximation for **equivalence relations** over access paths

Examples of aliasing abstractions:

- **set abstractions**: map from access paths to their equivalence class (ex: $\{\{p_0, p_1, \&x\}, \{p_2, p_3\}, \dots\}$)
- **numerical relations**, to describe aliasing among paths of the form $x(->n)^k$ (ex: $\{\{x(->n)^k, \&(x(->n)^{k+1}) \mid k \in \mathbb{N}\}$)

Weak update problems

```
x ∈ [-10, -5]; y ∈ [5, 10]
```

```
int * p;
```

```
if(?)
```

```
    p = &x;
```

```
else
```

```
    p = &y;
```

```
*p = 0;
```

- What is the final range for x ?
- What is the final range for y ?

Weak update problems

```
x ∈ [-10, -5]; y ∈ [5, 10]
int * p;
if(?)
    p = &x;
else
    p = &y;
*p = 0;
```

- What is the final range for x ?
- What is the final range for y ?

- After the **if** statement, p may contain any address in {&x, &y}
- Thus, the assignment must consider all cases, in a conservative way
- Thus, x may receive a new value (0) or keep its old value
- Conclusion: $x \in [-10, 0]$, $y \in [0, 10]$

Weak updates

Any imprecision in the analysis may lead to weak updates...

Limitation of basic pointer analyses

- **Weak updates:**
imprecisions for pointer values quickly spread out
- Many programs with pointers address **unbounded memory**
e.g., to create lists, trees and other dynamically allocated structures
most pointer analyses do not deal with this well...
- Pointer analyses do not nicely capture **structural invariants**
e.g., lists, trees, but also nested structures

Outline

- 1 Towards memory properties
- 2 Memory models
- 3 Abstraction of arrays
- 4 Abstraction of strings and buffers
- 5 Basic pointer analyses
- 6 Three valued logic heap abstraction**
 - Basic principles
 - Building an abstract domain
 - Weakening abstract elements
 - Computation of transfer functions

An abstract representation of memory states: shape graphs

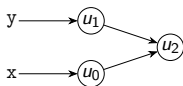
Goal of the static analysis

Discover complex invariants of programs that manipulate unbounded heap

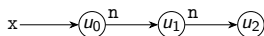
Observation: representation of memory states by shape graphs

- **Nodes** (aka, atoms) denote **memory locations**
- **Edges** denote **properties**, such as:
 - ▶ “field f of location u points to v ”
 - ▶ “variable x is stored at location u ”

Two alias pointers:



A list of length 2:



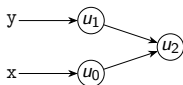
⇒ We need to over-approximate sets of shape graphs

Shape graphs and their representation

Description with predicates

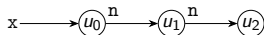
- **Boolean encoding:** nodes are atoms u_0, u_1, \dots
- **Predicates over atoms:**
 - ▶ $x(u)$: variable x contains the address of u
 - ▶ $n(u, v)$: field of u points to v
- **Truth values:** traditionally noted 0 and 1 in the TVLA litterature

Two alias pointers:



	x	y	\mapsto	u_0	u_1	u_2
u_0	1	0	u_0	0	0	1
u_1	0	1	u_1	0	0	1
u_2	0	0	u_2	0	0	0

A list of length 2:

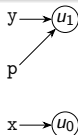
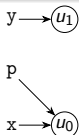


	x	$\cdot n \mapsto$	u_0	u_1	u_2
u_0	1	u_0	0	1	0
u_1	0	u_1	0	0	1
u_2	0	u_2	0	0	0

Unknown value: three valued logic

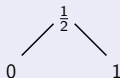
How to abstract away some information ?
i.e., to abstract several graphs into one ?

Example: pointer variable p
alias with x or y

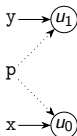


A boolean lattice

- Use **predicate tables**
- Add a \top boolean value;
(denoted to by $\frac{1}{2}$ in TVLA papers)



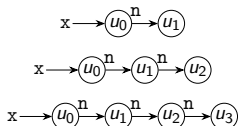
- Graph representation:
dotted edges
- **Abstract graph:**



Summary nodes

We cannot talk about unbounded memory states with finitely many nodes

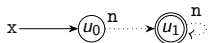
Lists of lengths 1, 2, 3:



We would like to **summarize** the lists

An idea

- Choose a node to represent **several** concrete nodes
- Similar to **smashing**



- Edges to u_1 are dotted

Definition: summary node

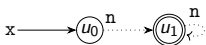
A **summary node** is an atom that may denote several concrete atoms

A few interesting predicates

We have already seen:

- $x(u)$: variable x contains the address of u
- $n(u, v)$: field of u points to v
- $\text{sum}(u)$: whether u is a summary node (convention: either 0 or $\frac{1}{2}$)

The properties of lists are not well-captured in



“Is shared”

$\text{sh}(u)$ ssi:

$$\exists v_0, v_1, \begin{cases} v_0 \neq v_1 \\ \wedge n(v_0, u) \\ \wedge n(v_1, u) \end{cases}$$

Predicates defined by transitive closure

- **Reachability**: $\underline{r}(u, v)$ ssi

$$u = v \vee \exists u_0, n(u, u_0) \wedge \underline{r}(u_0, v)$$

- **Acyclicity**: $\underline{\text{acy}}(v)$

similar, with a negation

Outline

- 1 Towards memory properties
- 2 Memory models
- 3 Abstraction of arrays
- 4 Abstraction of strings and buffers
- 5 Basic pointer analyses
- 6 Three valued logic heap abstraction**
 - Basic principles
 - Building an abstract domain
 - Weakening abstract elements
 - Computation of transfer functions

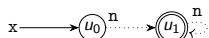
Three structures

Definition: 3-structures

A 3-structure is a tuple $(\mathcal{U}, \mathcal{P}, \phi)$:

- a set $\mathcal{U} = \{u_0, u_1, \dots, p_m\}$ of **atoms**
- a set $\mathcal{P} = \{p_0, p_1, \dots, p_n\}$ of **predicates**
(we write k_i for the arity of predicate p_i)
- a **truth table** ϕ such that $\phi(p_i, u_{l_1}, \dots, u_{l_{k_i}})$ denotes the truth value of p_i for $u_{l_1}, \dots, u_{l_{k_i}}$

note: truth values are elements of the lattice $\{0, \frac{1}{2}, 1\}$



$$\left\{ \begin{array}{l} \mathcal{U} = \{u_0, u_1\} \\ \mathcal{P} = \{x(\cdot), n(\cdot, \cdot), \underline{\text{sum}}(\cdot)\} \end{array} \right.$$

	x	<u>sum</u>
u_0	1	0
u_1	0	$\frac{1}{2}$
n	u_0	u_1
u_0	0	1
u_1	0	0

Embedding

- How to compare two 3-structures ?
- How to describe the concretization of 3-structures ?

The embedding principle

Let $\mathcal{S}_0 = (\mathcal{U}_0, \mathcal{P}, \phi_0)$ and $\mathcal{S}_1 = (\mathcal{U}_1, \mathcal{P}, \phi_1)$ be two three structures, with the same sets of predicates.

Let $f : \mathcal{U}_0 \rightarrow \mathcal{U}_1$, surjective.

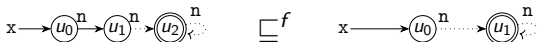
We say that f **embeds** \mathcal{S}_0 **into** \mathcal{S}_1 iff

$$\begin{aligned} &\text{for all predicate } p \in \mathcal{P} \text{ or arity } k, \\ &\quad \text{for all } u_{l_1}, \dots, u_{l_{k_i}} \in \mathcal{U}_0, \\ &\quad \phi_0(u_{l_1}, \dots, u_{l_{k_i}}) \sqsubseteq \phi_1(f(u_{l_1}), \dots, f(u_{l_{k_i}})) \end{aligned}$$

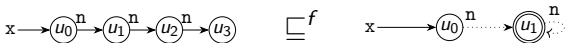
Then, **we write** $\mathcal{S}_0 \sqsubseteq^f \mathcal{S}_1$

Note: we use the order \sqsubseteq of the lattice $\{0, \frac{1}{2}, 1\}$

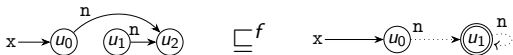
Embedding examples



where $f : u_0 \mapsto u_0; u_1 \mapsto u_1; u_2 \mapsto u_1$



where $f : u_0 \mapsto u_0; u_1 \mapsto u_1; u_2 \mapsto u_1; u_3 \mapsto u_1$



where $f : u_0 \mapsto u_0; u_1 \mapsto u_1; u_2 \mapsto u_1$

Note on the last example

- Reachability would be necessary to constrain it be a list
- Alternatively: cells should not be shared

Two structures and concretization

Concrete states correspond to 2-structures

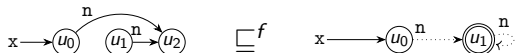
- **2-structure**: a 3-structure $(\mathcal{U}, \mathcal{P}, \phi)$ is a 2-structure, if and only if ϕ always returns in $\{0, 1\}$
- A **2-structure** corresponds to a set of concrete memory states (environment, heap):
 - ▶ we simply need to take into account all mappings of addresses into the memory
we let **stores**(\mathcal{S}) denote the stores corresponding to 2-structure \mathcal{S}
 - ▶ more on this in the next lecture; here we keep it informal

Concretization

$$\gamma(\mathcal{S}) = \bigcup \{ \text{stores}(\mathcal{S}') \mid \mathcal{S}' \text{ 2-structure s.t. } \exists f, \mathcal{S}' \sqsubseteq^f \mathcal{S} \}$$

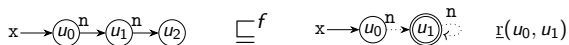
Concretization examples

- Without reachability:



where $f : u_0 \mapsto u_0; u_1 \mapsto u_1; u_2 \mapsto u_1; u_3 \mapsto u_1$

- With reachability:



where $f : u_0 \mapsto u_0; u_1 \mapsto u_1; u_2 \mapsto u_1$

Principle for the design of sound transfer functions

How to carry out static analysis using 3-structures ?

Embedding theorem

- Let $\mathcal{S}_0 = (\mathcal{U}_0, \mathcal{P}, \phi_0)$ and $\mathcal{S}_1 = (\mathcal{U}_1, \mathcal{P}, \phi_1)$ be two three structures, with the same sets of predicates
- Let $f : \mathcal{U}_0 \rightarrow \mathcal{U}_1$, such that $\mathcal{S}_0 \sqsubseteq^f \mathcal{S}_1$
- Let Ψ be a logical formula, with variables in X and $g : X \rightarrow \mathcal{U}_0$ be an assignment for the variables of Ψ

Then,
$$\llbracket \Psi \rrbracket_{|g}(\mathcal{S}_0) \sqsubseteq \llbracket \Psi \rrbracket_{|f \circ g}(\mathcal{S}_1)$$

Principle for the design of sound transfer functions

Transfer functions for static analysis

- Semantics of concrete statements encoded into boolean formulas
- **Example:** assignment $y := x$
 - ▶ let y' denote the *new* value of y
 - ▶ add the constraint $y'(u) = x(u)$
 - ▶ rename y' into y

Full examples of transfer functions computation in a few slides...

- Evaluation in the abstract is sound (embedding theorem)

Advantages:

- **abstract transfer functions** derive directly from the concrete transfer functions
intuition: $\alpha \circ f \circ \gamma \dots$
- the same solution works for **weakest pre-conditions**

Outline

- 1 Towards memory properties
- 2 Memory models
- 3 Abstraction of arrays
- 4 Abstraction of strings and buffers
- 5 Basic pointer analyses
- 6 Three valued logic heap abstraction**
 - Basic principles
 - Building an abstract domain
 - Weakening abstract elements
 - Computation of transfer functions

A powerset abstraction

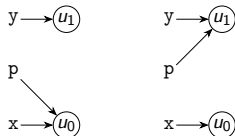
- Do 3-structures allow for a sufficient level of precision ?
- How to over-approximate a set of two-structures ?

```

int * x; f * y; ...
int * p = NULL;
if(...){
    p = x;
}else{
    p = y;
}
printf("%d", *p);
*xp = ...;

```

After the if statement:



abstracting here would be imprecise

Powerset abstraction

- Shape analyzers usually rely on a **powerset abstract domain** i.e., TVLA manipulates **finite disjunctions** of 3-structures
- How to ensure disjunctions will not grow infinite ?

Canonical abstraction

Canonicalization principle

Let \mathcal{L} be a lattice, $\mathcal{L}' \subseteq \mathcal{L}$ be a finite sub-lattice and $\mathbf{can} : \mathcal{L} \rightarrow \mathcal{L}'$:

- \mathbf{can} called a **canonicalization** if it is an upper closure operator
- then, \mathbf{can} extends into a canonicalization operator of $\mathcal{P}(\mathcal{L})$, into $\mathcal{P}(\mathcal{L}')$:

$$\mathbf{can}(\mathcal{E}) = \{\mathbf{can}(x) \mid x \in \mathcal{E}\}$$

To make the powerset domain work, we simply need a \mathbf{can} over 3-structures

A canonicalization over 3-structures

- We assume there are n variables x_1, \dots, x_n
Thus the number of unary predicates is finite
- **Sub-lattice**: structures with atoms **distinguished by the values of the unary predicates** (or *abstraction predicates*) x_1, \dots, x_n

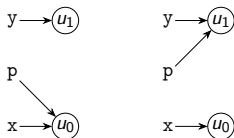
We may choose another set of predicates for the sub-lattice representation

Canonical abstraction

Canonical abstraction by truth blurring

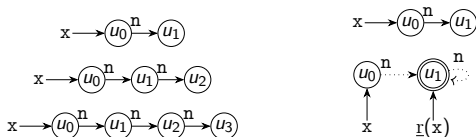
- 1 **Identify** nodes that **have different abstraction predicates**
- 2 When several nodes have the **same abstraction predicate** introduce a **summary node**
- 3 **Compute new predicate values** by doing a join over truth values

Elements not merged:



Elements merged:

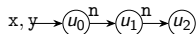
Lists of lengths 1, 2, 3: **Abstract into:**



Outline

- 1 Towards memory properties
- 2 Memory models
- 3 Abstraction of arrays
- 4 Abstraction of strings and buffers
- 5 Basic pointer analyses
- 6 Three valued logic heap abstraction**
 - Basic principles
 - Building an abstract domain
 - Weakening abstract elements
 - Computation of transfer functions

Assignment: a simple case

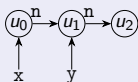
Statement $\ell_0 : y = y \rightarrow n; \ell_1 : \dots$ Pre-condition \mathcal{S} 

Transfer function

- Should yield an over-approximation of $\{m_1 \in \mathbb{M} \mid (\ell_0, m_0) \rightarrow (\ell_1, m_1)\}$
- We let **“primed predicates”** denote predicates after evaluation of the assignment, to evaluate them in the same structure
- Then:

$$\begin{aligned} x'(u) &= x(u) \\ y'(u) &= \exists v, y(v) \wedge n(v, u) \\ n'(u, v) &= n(u, v) \end{aligned}$$

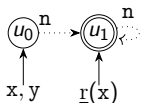
- Result:



This was exactly what we expected

Assignment: a more involved case

Statement $l_0 : y = y \rightarrow n; l_1 : \dots$



Pre-condition \mathcal{S}

- Let us try to resolve the update in the same way as before:

$$\begin{aligned} x'(u) &= x(u) \\ y'(u) &= \exists v, y(v) \wedge n(v, u) \\ n'(u, v) &= n(u, v) \end{aligned}$$

- We **cannot resolve** y' :

$$\begin{cases} y'(u_0) = 0 \\ y'(u_1) = \frac{1}{2} \end{cases}$$

Imprecision: after the statement, y may point to anywhere in the list, save for the first element...

- The assignment transfer function cannot be computed immediately
- We need to refine the 3-structure first

Focus

Focusing on a formula

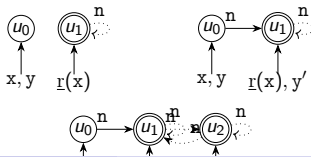
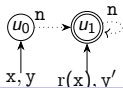
We assume a 3-structure \mathcal{S} and a boolean formula f are given, we call a **focusing** \mathcal{S} on f the generation of a set $\hat{\mathcal{S}}$ such that:

- f evaluates to 0 or 1 on all elements of $\hat{\mathcal{S}}$
- **precision was gained:** $\forall \mathcal{S}' \in \hat{\mathcal{S}}, \mathcal{S}' \sqsubseteq \mathcal{S}$
- **soundness is preserved:** $\gamma(\mathcal{S}) = \bigcup \{ \gamma(\mathcal{S}') \mid \mathcal{S}' \in \hat{\mathcal{S}} \}$

- Focusing algorithms are complex and tricky (see biblio)
- Involves splitting of summary nodes, solving of boolean constraints

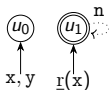
We obtain (we show y and y'):

Example: focusing on
 $y'(u) = \exists v, y(v) \wedge n(v, u)$

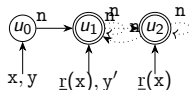


Focus and coerce

Some of the 3-structures generated by focus are not precise



u_1 is reachable from x , but there is no sequence of n fields: this structure has **empty concretization**

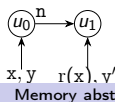


u_0 has an n -field to u_1 so u_1 denotes a unique atom and **cannot be a summary node**

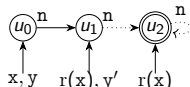
Coerce operation

It enforces logical constraints among predicates and discards 3-structures with an empty concretization

Result:



Memory abstraction



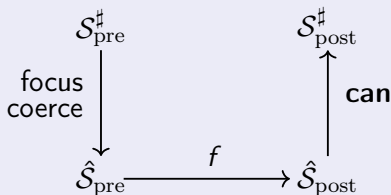
Dec, 10th. 2014

92 / 96

Focus, transfer, abstract...

Computation of a transfer function

We consider a transfer function encoded into boolean formula f



Soundness proof steps:

- ① sound encoding of the semantics of program statements into formulas typically, no loss of precision at this stage
- ② focusing should yield an over-approximation of its input
- ③ canonicalization over-approximates graph (truth blurring weakening)

A common picture in shape analysis

Outline

- 1 Towards memory properties
- 2 Memory models
- 3 Abstraction of arrays
- 4 Abstraction of strings and buffers
- 5 Basic pointer analyses
- 6 Three valued logic heap abstraction
- 7 Conclusion**

Conclusion

Concrete semantics:

- Splitting environment and heap
- Taking into account of the representaton of data

Many families of domain specific abstractions:

- Based on numerical methods
path based pointer analyses, array segment analyses, string analyses
- Symbolic abstractions based on pointer sets, structural predicates
- Locally concretize / globally abstract pattern (TVLA, arrays...)
More on this during the next lecture...

Bibliography

- **[HP'08]: Discovering properties about arrays in simple programs.**
Nicolas Halbwachs, Mathias Péron. In PLDI'08, pages 339-348, 2008.
- **[CCL'11]: A parametric segmentation functor for fully automatic and scalable array content analysis.**
Patrick Cousot, Radhia Cousot, Francesco Logozzo. In POPL'11, pages 105-118, 2011.
- **[AD'94]: Interprocedural may alias analysis for pointers: beyond k -limiting.**
Alain Deutsch. In PLDI'94, pages 230–241, 1994.
- **[CSSV'03]: CSSV: towards a realistic tool for statically detecting all buffer overflows in C.**
Nurit Dor, Michael Rodeh, Shmuel Sagiv. In PLDI'03, pages 155-167.
- **[SRW'99]: Parametric Shape Analysis via 3-Valued Logic.**
Shmuel Sagiv, Thomas W. Reps et Reinhard Wilhelm. In POPL'99, pages 105–118, 1999.