

# Design of a Modular Platform for Static Analysis

Antoine Miné    Abdelraouf Ouadjaout\*    Matthieu Journault  
Raphaël Monat

LIP6, Sorbonne Université, Paris

Journée de l'équipe APR  
July 5, 2019

# What is MOPSA?

MOPSA is a **Modular Open Platform for Static Analysis**

## 1 Based on **abstract interpretation**

- ∇ **Sound** Cover all possible executions.
- 🕒 **Efficient** Terminates in finite time.
- 💻 **Automatic** No need to manual annotation.

## 2 Supports **multiple languages**

- 📄 **Expressiveness** Keep the original AST of the program.
- ♻️ **Reusability** Reuse abstractions among languages.

## 3 Features a **flexible architecture**

- 🧩 **Loose coupling** Divided into interchangeable components.
- 🧱 **Composition** Create complex components from simpler ones.
- 💬 **Cooperation** Components can communicate and delegate tasks.

🗨️ May raise false alarms.

# Abstract Interpretation

How can an abstract interpreter handle this program?

```
int main(int argc, char *argv[]) {  
    char **p = argv; int i = 0;  
    while(*p) {  
        printf("length(argv[%d]) = %lu\n", i, strlen(*p));  
        p++; i++;  
    }  
}
```

# Abstract Interpretation

How can an abstract interpreter handle this program?

```
● int main(int argc, char *argv[]) {  
    char **p = argv; int i = 0;  
    while(*p) {  
        printf("length(argv[%d]) = %lu\n", i, strlen(*p));  
        p++; i++;  
    }  
}
```

- Use computable abstractions to over-approximate states.

Numeric

$\text{argc} \in [1, \text{SIZE\_MAX} - 1]$

$\text{size}(\text{argv}) = \text{argc} + 1$

# Abstract Interpretation

How can an abstract interpreter handle this program?

```
● int main(int argc, char *argv[]) {  
    char **p = argv; int i = 0;  
    while(*p) {  
        printf("length(argv[%d]) = %lu\n", i, strlen(*p));  
        p++; i++;  
    }  
}
```

- Use computable abstractions to over-approximate states.

Numeric

$\text{argc} \in [1, \text{SIZE\_MAX} - 1]$

$\text{size}(\text{argv}) = \text{argc} + 1$

$\text{size}(\text{⓪}) \in [1, \text{SIZE\_MAX}]$

Pointers

$\text{argv}[0..\text{argc} - 1] \rightsquigarrow \{\text{⓪}\}$

$\text{argv}[\text{argc}] \rightsquigarrow \text{NULL}$

# Abstract Interpretation

How can an abstract interpreter handle this program?

```
● int main(int argc, char *argv[]) {  
    char **p = argv; int i = 0;  
    while(*p) {  
        printf("length(argv[%d]) = %lu\n", i, strlen(*p));  
        p++; i++;  
    }  
}
```

- Use computable abstractions to over-approximate states.

| Numeric                                           | Pointers                                                        | Strings                                                          |
|---------------------------------------------------|-----------------------------------------------------------------|------------------------------------------------------------------|
| $\text{argc} \in [1, \text{SIZE\_MAX} - 1]$       | $\text{argv}[0..\text{argc} - 1] \rightsquigarrow \{\text{@}\}$ | $\exists k \in [0, \text{size}(\text{@}) - 1] : \text{@}[k] = 0$ |
| $\text{size}(\text{argv}) = \text{argc} + 1$      | $\text{argv}[\text{argc}] \rightsquigarrow \text{NULL}$         |                                                                  |
| $\text{size}(\text{@}) \in [1, \text{SIZE\_MAX}]$ |                                                                 |                                                                  |

# Abstract Interpretation

How can an abstract interpreter handle this program?

```
int main(int argc, char *argv[]) {
    char **p = argv; int i = 0;
    while(*p) {
        printf("length(argv[%d]) = %lu\n", i, strlen(*p));
        p++; i++;
    }
}
```

- Use computable abstractions to over-approximate states.
- Infer inductive invariants.

| Numeric                                                     | Pointers                                                        | Strings                                                          |
|-------------------------------------------------------------|-----------------------------------------------------------------|------------------------------------------------------------------|
| $\text{argc} \in [1, \text{SIZE\_MAX} - 1]$                 | $\text{argv}[0..\text{argc} - 1] \rightsquigarrow \{\text{@}\}$ | $\exists k \in [0, \text{size}(\text{@}) - 1] : \text{@}[k] = 0$ |
| $\text{size}(\text{argv}) = \text{argc} + 1$                | $\text{argv}[\text{argc}] \rightsquigarrow \text{NULL}$         |                                                                  |
| $\text{size}(\text{@}) \in [1, \text{SIZE\_MAX}]$           | $p \rightsquigarrow \{\text{argv}\}$                            |                                                                  |
| $0 \leq \text{offset}(p) \leq \text{size}(\text{argv}) - 1$ |                                                                 |                                                                  |
| $\text{offset}(p) = i$                                      |                                                                 |                                                                  |

# Abstract Interpretation

How can an abstract interpreter handle this program?

```
int main(int argc, char *argv[]) {
    char **p = argv; int i = 0;
    while(*p) {
        printf("length(argv[%d]) = %lu\n", i, strlen(*p));
        p++; i++;
    }
}
```

Valid string passed to strlen?

Integer overflow on i?

- Use computable abstractions to over-approximate states.
- Infer inductive invariants.
- Check safety rules.

| Numeric                                              | Pointers                                           | Strings                                                            |
|------------------------------------------------------|----------------------------------------------------|--------------------------------------------------------------------|
| $argc \in [1, \text{SIZE\_MAX} - 1]$                 | $argv[0..argc - 1] \rightsquigarrow \{\emptyset\}$ | $\exists k \in [0, \text{size}(\emptyset) - 1] : \emptyset[k] = 0$ |
| $\text{size}(argv) = argc + 1$                       | $argv[argc] \rightsquigarrow \text{NULL}$          |                                                                    |
| $\text{size}(\emptyset) \in [1, \text{SIZE\_MAX}]$   | $p \rightsquigarrow \{argv\}$                      |                                                                    |
| $0 \leq \text{offset}(p) \leq \text{size}(argv) - 1$ |                                                    |                                                                    |
| $\text{offset}(p) = i$                               |                                                    |                                                                    |



## Part I

# Multi-Language Support

# Multi-Language Support

## Extensible AST

### Existing Solutions

- **Static** translation to an intermediate language.  
*e.g. LLVM, SIL of Facebook Infer(C/C++/ObjectiveC/Java).*
- 👍 Easy reuse of abstractions.
- 👎 Loss of source code information.

### MOPSA Approach

- Keep the original AST of the program.
- Domains can **extend** the language to add new AST nodes.

```
type expr_kind += | E_c_call of ...           type expr_kind += | E_py_call of ...
```

- Intermediate languages are added to share abstractions.  
*Original and intermediate languages are in the same AST.*
- Translation is **dynamic**, **semantic** and **cooperative**.

# Multi-Language Support

## Cooperation via Delegation

```
int f(int x)
{
  <...>
}
```

```
class F:
    def __call__(self, x):
        <...>
f = F()
```

man.eval (E\_c\_call(f, [x]))

man.eval (E\_py\_call(f, [x]))

### C.Interproc

- Evaluate f (in case of function pointer)
- Get body and check args
- man.eval (E\_u\_call(body, [x]))

### Python.Interproc

- Evaluate the object f
- Check that C defines \_\_call\_\_
- Get body and check args
- man.eval (E\_u\_call(body, [f; x]))

### Universal.Interproc

- Match E\_u\_call(body, args)
- Use inlining, summaries, etc.

## Part II

# Flexible Architecture

# Flexible Architecture

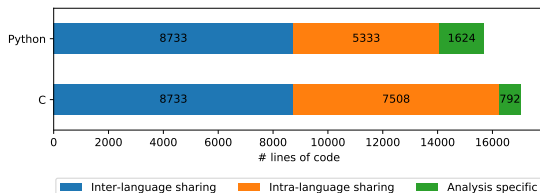
## Domains

### Domains

- Analyzer is decomposed into **domains** over extensible AST.
- Share domains as much as possible.
- Several domains can abstract the same thing, but differently.  
 $\{\langle x \mapsto 0 \rangle, \langle x \mapsto 2 \rangle\}$  can be abstracted as  $\langle x \mapsto [0, 2] \rangle$  or  $\langle x \mapsto 2\mathbb{Z} \rangle$ .

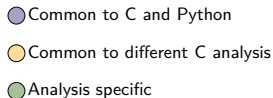
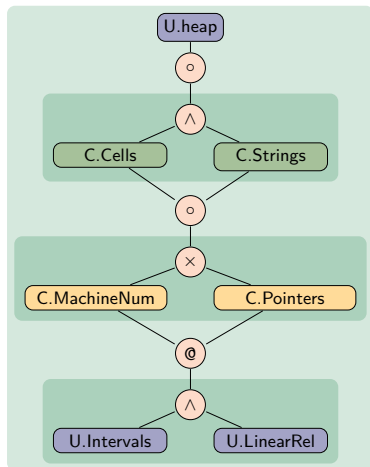
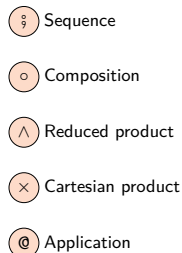
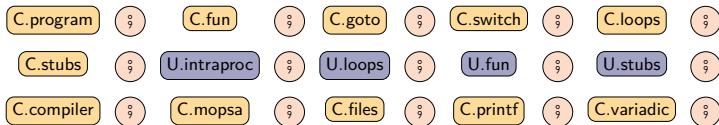
### Unified Signature

- MOPSA defines a **unified** signature for all domains.
- Other simplified signatures are automatically lifted.



# Flexible Architecture

## Composition

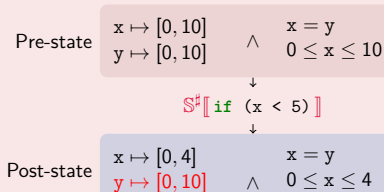


INTERVALS  $\stackrel{\text{def}}{=} \mathcal{V} \rightarrow (\mathbb{Z} \cup \{-\infty\} \times \mathbb{Z} \cup \{+\infty\})$

LINEARREL  $\stackrel{\text{def}}{=} \bigwedge_j (\sum_{i=1}^n \alpha_{ij} v_i \geq \beta_j), \alpha_{ij}, \beta_j \in \mathbb{R}, v_i \in \mathcal{V}$

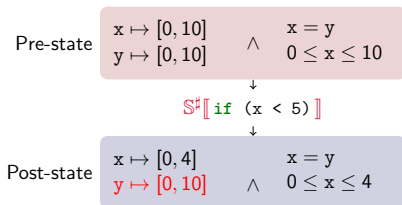
- INTERVALS is a fast domain provided by most analyzers.
- LINEARREL keeps track of linear relations.  
e.g. *size*(argv) = argc + 1
- LINEARREL and INTERVALS have the same concrete semantics.  
 $\implies$  can be combined with a meet product  $\wedge$ .

## Need for Collaboration



- INTERVALS can not infer that  $y < 5$ .
- How to collaborate while preserving loose coupling?

- Reduction rules combine the results of two abstract domains.
- They are implemented **outside** domains to reduce coupling.



## The hidden constraint reduction

- 1 Collect variables  $V$  statically present in the statement.  
 $V = \{x\}$
- 2 Search for the set variables  $R$  in LINEARREL related to  $V$ .  
 $R = \{y\}$ , since  $x = y$
- 3 Get the interval of  $R$  in LINEARREL and refine INTERVALS.  
 $\text{INTERVALS}[x \mapsto [0, 10] \sqcap [0, 4]]$



# Scalar Domains

## Machine Numbers

MachineNum @ (LinearRel ^ Intervals)

- INTERVALS and LINEARREL uses **mathematical** numbers.
- C uses **finite precision** numbers with modular arithmetics.
- MACHINENUM lifts C statements to a math semantic.  
→ Operator @ applies a **stack** domain to an argument.

# Scalar Domains

## Machine Numbers

MachineNum   @   (LinearRel   ^   Intervals)

- INTERVALS and LINEARREL uses **mathematical** numbers.
- C uses **finite precision** numbers with modular arithmetics.
- MACHINENUM lifts C statements to a math semantic.  
→ Operator @ applies a **stack** domain to an argument.

## Dynamic Lifting

Consider a variable  $x$  declared as **unsigned char**.

$$x + 2 \in [0, 255] \rightarrow \{x + 2, \langle x \mapsto [0, 253] \rangle, \emptyset\}$$

$$\mathbb{E}^\# \llbracket x + 2 \rrbracket \langle x \mapsto [0, 255] \rangle$$

$$x + 2 \notin [0, 255] \rightarrow \{wrap(x + 2, [0, 255]), \langle x \mapsto [254, 255] \rangle, \{IntOverflow\}\}$$

# Scalar Domains

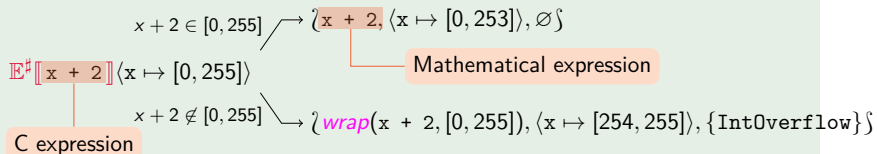
## Machine Numbers

MachineNum   @   (LinearRel   ^   Intervals)

- INTERVALS and LINEARREL uses **mathematical** numbers.
- C uses **finite precision** numbers with modular arithmetics.
- MACHINENUM lifts C statements to a math semantic.  
→ Operator @ applies a **stack** domain to an argument.

## Dynamic Lifting

Consider a variable  $x$  declared as **unsigned char**.



### Result Partitioning

MOPSA uses **monads and delegation** to handle  $\{r_1\} \vee \dots \vee \{r_n\}$ .

```
let eval exp man a =
  match exp.ekind with
  | E_var v -> Some (Eval.singleton exp a)
  | E_binop(op, e1, e2) ->
    man.eval e1 a >>= fun n1 a ->
    man.eval e2 a >>= fun n2 a ->
    let vmin, vmax = rangeof exp.ety in
    let nexp = mk_binop n1 op n2 in
    let ret = assume (mk_in nexp vmin vmax) man a
      ~fthen:(fun a -> Eval.singleton nexp a)
      ~felse:(fun a ->
        let nexp' = mk_wrap nexp vmin vmax in
        raise_alarm IntegerOverflow a |> Eval.singleton nexp'
      )
    in
    Some ret
  | _ -> None
```

### Result Partitioning

MOPSA uses **monads and delegation** to handle  $\{r_1\} \vee \dots \vee \{r_n\}$ .

```

let eval exp man a =
  match exp.ekind with
  | E_var v -> Some (Eval.singleton exp a) — Variables do not overflow
  | E_binop(op, e1, e2) ->
    man.eval e1 a >>= fun n1 a ->
    man.eval e2 a >>= fun n2 a ->
    let vmin, vmax = rangeof exp.ety in
    let nexp = mk_binop n1 op n2 in
    let ret = assume (mk_in nexp vmin vmax) man a
      ~fthen:(fun a -> Eval.singleton nexp a)
      ~felse:(fun a ->
        let nexp' = mk_wrap nexp vmin vmax in
        raise_alarm IntegerOverflow a |> Eval.singleton nexp'
      )
    in
    Some ret
  | _ -> None
  
```

### Result Partitioning

MOPSA uses **monads and delegation** to handle  $\{r_1\} \vee \dots \vee \{r_n\}$ .

```

let eval exp man a =
  match exp.ekind with
  | E_var v -> Some (Eval.singleton exp a)
  | E_binop(op, e1, e2) ->
    man.eval e1 a >>= fun n1 a -> — Evaluate e1 and bind each case {n1}
    man.eval e2 a >>= fun n2 a ->
      let vmin, vmax = rangeof exp.etyt in
      let nexp = mk_binop n1 op n2 in
      let ret = assume (mk_in nexp vmin vmax) man a
        ~fthen:(fun a -> Eval.singleton nexp a)
        ~felse:(fun a ->
          let nexp' = mk_wrap nexp vmin vmax in
          raise_alarm IntegerOverflow a |> Eval.singleton nexp'
        )
      in
      Some ret
  | _ -> None
  
```

### Result Partitioning

MOPSA uses **monads and delegation** to handle  $\{r_1\} \vee \dots \vee \{r_n\}$ .

```

let eval exp man a =
  match exp.ekind with
  | E_var v -> Some (Eval.singleton exp a)
  | E_binop(op, e1, e2) ->
    man.eval e1 a >>= fun n1 a ->
      man.eval e2 a >>= fun n2 a -> — Evaluate e2 and bind each case {n2}
      let vmin, vmax = rangeof exp.ety in
      let nexp = mk_binop n1 op n2 in
      let ret = assume (mk_in nexp vmin vmax) man a
        ~fthen:(fun a -> Eval.singleton nexp a)
        ~felse:(fun a ->
          let nexp' = mk_wrap nexp vmin vmax in
          raise_alarm IntegerOverflow a |> Eval.singleton nexp'
        )
      in
      Some ret
  | _ -> None
  
```

### Result Partitioning

MOPSA uses **monads and delegation** to handle  $\{r_1\} \vee \dots \vee \{r_n\}$ .

```

let eval exp man a =
  match exp.ekind with
  | E_var v -> Some (Eval.singleton exp a)
  | E_binop(op, e1, e2) ->
    man.eval e1 a >>= fun n1 a ->
    man.eval e2 a >>= fun n2 a ->
    let vmin, vmax = rangeof exp.ety in
    let nexp = mk_binop n1 op n2 in
    let ret = assume (mk_in nexp vmin vmax) man a
      ~fthen:(fun a -> Eval.singleton nexp a)
      ~felse:(fun a ->
        let nexp' = mk_wrap nexp vmin vmax in
        raise_alarm IntegerOverflow a |> Eval.singleton nexp'
      )
    in
    Some ret
  | _ -> None
  
```

Partition on condition  
 $nexp \in [vmin, vmax]$





$\text{POINTERS}(\text{NUM}) \stackrel{\text{def}}{=} (\mathcal{V}_{ptr} \rightarrow \wp(\mathcal{V}) \cup \{\text{NULL}, \text{INVALID}\}) \times \text{NUM}$

- POINTERS is also a **stack** domain.

- 1 Each pointer is mapped to the set of pointed variables.
- 2 The offset of a pointer  $p$  is abstracted by *offset*( $p$ ) in the argument numeric domain NUM.

# Scalar Domains

## Pointers

(Pointer) × (MachineNum) @ (LinearRel) ∧ (Intervals)

$\text{POINTERS}(\text{NUM}) \stackrel{\text{def}}{=} (\mathcal{V}_{\text{ptr}} \rightarrow \wp(\mathcal{V}) \cup \{\text{NULL}, \text{INVALID}\}) \times \text{NUM}$

- POINTERS is also a **stack** domain.
  - 1 Each pointer is mapped to the set of pointed variables.
  - 2 The offset of a pointer  $p$  is abstracted by *offset*( $p$ ) in the argument numeric domain NUM.
- POINTERS and MACHINE\_NUM lift a **shared** numeric environment using a cartesian product (×).
  - 👍 Infer relations between offsets and numeric variables.

### Example

```
char a[10] = "hello";  
int i = _mopsa_rand(0,9);  
char *p = &(a[i]); /* ⟨p ↦ {a}⟩, ⟨i ∈ [0,9] ∧ offset(p) = i⟩ */
```

# Block Domains



## CELLS

- Represent each scalar sub-block individually.
- Tailored for low-level C (casts, type punning, unions).

```
union { uint16 ax; struct { uint8 al; uint8 ah; } bytes; } reg;  
reg.ax = 0xabcd; /* <reg[0 : 2] = 43981> */  
x = reg.bytes.al; /* <x = 205> */
```

## STRINGS

- Keep track of the first ' $\backslash 0$ ' byte in the block.
- Useful for checking valid strings.

```
char s[100];  
int i, n = _mopsa_rand(1,10);  
for(i = 0; i < n; i++) s[i] = 'a';  
s[i] = '\0'; /* <n  $\in$  [1, 10]  $\wedge$  length(s) = n = i> */
```

# Model of the C Standard Library

- One goal of MOPSA is to analyze large open-source projects.
- C standard library is used in (almost) everywhere.
- MOPSA uses annotations inspired from Frama-C's ACSL.
  - 1 Specify the behavior of the APIs.
  - 2 Express safety functional properties.

```
/*$  
 * requires:  $\exists$  int i  $\in$  [0, size(__s) - 1]: __s[i] == 0;  
 * ensures : return  $\in$  [0, size(__s) - 1];  
 * ensures : __s[return] == 0;  
 * ensures : return > 0  $\implies$   
 *            $\forall$  int i  $\in$  [0, return - 1]: __s[i] != 0;  
 */  
size_t strlen(const char *__s);
```

# Experiments

## C Benchmarks

### Juliet test suite for Common Weakness Enumerations (CWE)

| Category | Description         | Lines | Time     | Alarms | Coverage |
|----------|---------------------|-------|----------|--------|----------|
| CWE476   | Null pointer deref. | 25k   | 2min26s  | 0      | 100%     |
| CWE369   | Division by zero    | 109k  | 7min20s  | 699    | 53%      |
| CWE190   | Integer overflow    | 440k  | 34min57s | 760    | 73%      |

## Python Benchmarks

### Official Python performance benchmarks

| Program       | Lines | Time  | Alarms |
|---------------|-------|-------|--------|
| fannkuch.py   | 59    | 0.07s | 0      |
| float.py      | 63    | 0.06s | 0      |
| spectral_n.py | 74    | 0.33s | 1      |
| nbody.py      | 157   | 1.5s  | 1      |
| chaos.py      | 324   | 5.9s  | 1      |

## Features

- MOPSA provides a compositional and flexible architecture for building sound static analyzers.
- Already used in research projects for the analysis on non-trivial C and Python programs.
- Support of different languages, reusable abstract domains with loose coupling.

## Limitations and Future Work

- Improve coverage of builtin functions and standard libraries.
- Not tested on large codebases.
- Open questions: backward analysis, shape analysis.