# Mopsa Tutorial

Raphaël Monat

`rmonat.fr/lipari24/`

Inria

Modular Open Platform for Static Analysis  [Jou+19]
gitlab.com/mopsa/mopsa-analyzer

Goals: explore new designs, ease development of (relational) analyses

### One AST to rule them all

- Multilanguage support
- Expressiveness
- Reusability

### Unified domain signature

- Semantic rewriting
- Loose coupling
- Observability

### DAG of abstractions

- Relational domains
- Composition
- Cooperation

# A motivating example

### Averaging tasks

```python
class Task:
  def __init__(self, weight):
    if weight < 0: raise ValueError
    self.weight = weight

def average(l):
  m = 0
  for i in range(len(l)):
    m = m + l[i].weight
  m = m // (i + 1)
  return m

l = [Task(randint(0, 20))
    for i in range(randint(5, 10))]
m = average(l)
```

Proved safe?

▶ `m // (i+1)`

▶ `l[i].weight`

Searching for a loop invariant
Stateless domains: list content, list length

### Environment abstraction

$m \mapsto @^{\sharp}_{\text{int}^{\sharp}}$   $i \mapsto @^{\sharp}_{\text{int}^{\sharp}}$   $\underline{\text{els}}(l) \mapsto @^{\sharp}_{\text{Task}}$
$\underline{@^{\sharp}_{\text{Task}} \cdot \texttt{weight}} \mapsto @^{\sharp}_{\text{int}^{\sharp}}$

### Numeric abstraction (polyhedra)

$m \in [0, +\infty)$   $\underline{\text{els}}(l) \in [0, 20]$
$0 \le i < \underline{\text{len}}(l)$   $5 \le \underline{\text{len}}(l) \le 10$
$0 \le \underline{@^{\sharp}_{\text{Task}} \cdot \texttt{weight}} \le 20$

### Attributes abstraction

$@^{\sharp}_{\text{Task}} \mapsto (\{\texttt{weight}\}, \emptyset)$

2

## Contributors (2018–2024, chronological arrival order)

- ▶ A. Miné
- ▶ A. Ouadjaout
- ▶ M. Journault
- ▶ A. Fromherz

- ▶ D. Delmas
- ▶ R. Monat
- ▶ G. Bau
- ▶ F. Parolini

- ▶ M. Milanese
- ▶ M. Valnet
- ▶ J. Boillot

Maintainers in bold.

Slides inspired from many contributors, mistakes are my own.

# Works around Mopsa

## Languages

C [JMO18; OM20], Python [MOM20a; MOM20b]
Multilanguage Python+C [MOM21]

WIP: Michelson [Bau+22], OCaml [VMM23], Catala (date arithmetic [MFM24])...

## Properties

▶ Absence of RTEs
▶ Patch analysis [DM19]
▶ Endianness portability [DOM21]
▶ Non-exploitability [PM24]
▶ Sufficient precondition inference [MM24]

## Software Verification Competition

We won the "SoftwareSystems" track of SV-Comp 2024 [Mon+24]!



5

Mopsa is implemented in OCaml.

▶ Curried functions: `f x y z` $\rightsquigarrow f(x, y, z)$

▶ `fun x -> e` $\Leftrightarrow \lambda x.e$

▶ Variable binding `let x = e1 in e2`

▶ Algebraic datatypes and pattern matching

```ocaml
type 'a option = None | Some of 'a

match e with
| None -> e1
| Some y -> e2 y
```

Polymorphism = Type Variables

6

# Outline

Then: 60 minutes practical session

## Objectives

► Understand Mopsa's core principles
► Ability to run C/Python analyses

**Feel free to ask questions!**

## New users welcome!

Feel free to reach out in the future!

7

# Architecture of Mopsa

Defining analyses

# Analysis = composition of abstract domains

unified domain signature $\implies$ iterators are abstract domains

▶ flexible architecture suitable for various programming paradigms

▶ separation of concerns

▶ allows reuse of domains across languages

▶ defined as json files in `share/mopsa/configs`



8

# Abstract state & domain signature

## Which type can we give to the abstract state?

- ▶ Polymorphism to the rescue: `'a` represents the abstract state
- ▶ Extended into `'a flow` to maintain additional info (more later)

## Handling cases

- ▶ `type ('a, 'r) cases` as DNFs over `'a flow * 'r`
- ▶ `Cases.singleton : 'r -> 'a flow -> ('a, 'r) cases`
- ▶ Binding operator `cases >>$ fun r flow -> ...`
  `>>$ : ('a,'r) cases -> ('r -> 'a flow -> ('a,'s) cases) -> ('a,'s) cases`
- ▶ Side note: this is a monad

# Abstract state & domain signature – II

The manager: interoperating the whole analysis and local domains

- ▶ Local domain has a private `type t`
- ▶ Whole abstract state of `type 'a`

$\implies$ `type ('a, t) man`

## From global analysis to local domain

- ▶ Get the domain's data
  `get : 'a -> t`
- ▶ Set the domain's data
  `set : t -> 'a -> 'a`

## From local domain to global analysis

- ▶ Analyze a given expression
- ▶ Analyze a given statement
  `man.exec stmt` $\sigma \Leftrightarrow \mathbb{S}^\sharp [\![ \text{stmt} ]\!] \sigma$

Signatures later

Also: lattice operators

# Abstract state & domain signature – III

## Utilities

```
1  type ('a, 'r) cases (* ≃ DNF of 'a flow * 'r *)
2
3  type 'a eval = ('a, expr) cases
4  type 'a post = ('a, unit) cases
5
6  (* Manager, allowing interaction between a
7     domain ('t) and whole analysis ('a) *)
8  type ('a, 't) man = {
9    get : 'a -> 't;
10   set : 't -> 'a -> 'a;
11   exec : stmt -> 'a flow -> 'a post;
12   eval : expr -> 'a flow -> 'a eval;
13   (* [...] *)
14 }
```

## Domain type overview

```
1  module type DOMAIN = sig
2
3  type t
4  (* private, opaque data of the domain *)
5  val name : string
6
7  val join : t -> t -> t (* and other lattice operators *)
8
9  (* Transfer functions *)
10 val exec : stmt -> ('a, t) man -> 'a flow -> 'a post option
11 val eval : expr -> ('a, t) man -> 'a flow -> 'a eval option
12
13 (* [...] *)
14 end
```

```
val exec : stmt -> ('a, t) man -> 'a flow -> 'a post option
val eval : expr -> ('a, t) man -> 'a flow -> 'a eval option
```

▶ `('a, t) man` manager

▶ `'a flow` abstract state

▶ option: domains return `None` for unsupported statements/expressions.

- `'a post = ('a, unit) cases`. DNF of abstract states.
- `'a eval = ('a, expr) cases`. DNF of abstract states and symbolic expressions. Useful for rewriting, esp. for relational analyses

## Example: loop iterator

Iterators are <u>stateless</u> domains:
- ▶ `type t = unit`, trivial lattice operators
- ▶ We have an automatic lifter for all that!

The loop iterator focuses on postfixpoint computation, and delegates the rest.

<div align="center">Universal.Iterators.Loops</div>

```
1  let rec lfp man cond body flow_init flow =
2    man.exec (mk_block [mk_assume cond; body]) flow >>$ fun () flow' ->
3    if man.lattice.subset (man.lattice.join flow_init flow') flow
4    then Cases.singleton () flow'
5    else lfp man cond body flow_init (man.lattice.widen flow flow')
6
7  let exec stmt man flow = match stmt.skind with
8    | S_while (cond, body) ->
9      Some (lfp man cond body flow flow >>$ fun () lfp_flow ->
10           man.exec (mk_assume (mk_not cond)) lfp_flow)
11   | _ -> None
```

# Architecture of Mopsa

Delegation-based support of multiple languages

# Iterators to handle multiple languages

## Traditional approaches

Desugar/compile programs to an intermediate representation (IR)
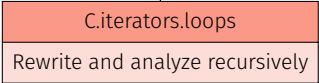
Example: Infer's IR has five (!) constructors

## Mopsa

▶ No loss of precision from the frontend[1]

▶ Various programming paradigms supported!

▶ All constructs have to be handled – but rewritings are possible

▶ A single AST type which can be extended for new languages

[1]By default, 3-address code may result in precision loss [NP18]

```
for(init; cond; incr) body
```

| C.iterators.loops |
|---|
| Rewrite and analyze recursively |

```
init;
while(cond) {
  body;
  incr;
}
```

```
for target in iterable: body
```

| Python.Desugar.Loops |
|---|
| o Rewrite and analyze recursively<br>o Optimize for some <u>semantic</u> cases |

```
it = iter(iterable)
while(1) {
 try: target = next(it)
 except StopIteration: break
 body
}
clean it
```

| Universal.Iterators.Loops |
|---|
| Matches while(...){...}<br>Computes fixpoint using widening |

15

# Architecture of Mopsa

Stateful domains through numerical examples

$$\mathcal{P}(\mathcal{V} \to \mathbb{Z}) \underset{\alpha}{\overset{\gamma}{\longleftrightarrow}} \underbrace{\mathcal{V} \to \mathcal{P}(\mathbb{Z})}_{\text{Cartesian abstraction}} \underset{\dot{\alpha}_l}{\overset{\dot{\gamma}_l}{\longleftrightarrow}} \underbrace{\mathcal{V} \to \mathcal{I}}_{\text{Lifting intervals}}$$

$$\alpha : \begin{cases} \mathcal{P}(\mathcal{V} \to \mathbb{Z}) & \to & \mathcal{V} \to \mathcal{P}(\mathbb{Z}) \\ \sum & \mapsto & \lambda v.\{\, \sigma(v) \mid \sigma \in \sum \,\} \end{cases} \qquad \gamma : \begin{cases} \mathcal{V} \to \mathcal{P}(\mathbb{Z}) & \to & \mathcal{P}(\mathcal{V} \to \mathbb{Z}) \\ f & \mapsto & \{\, \sigma \mid \forall v \in \mathcal{V}, \sigma(v) \in f(v) \,\} \end{cases}$$

$\implies$ We can define abstract operations on <u>values only</u>.

The `nonrel` combinator will lift values to the mapping.

16

# Numerical abstract values – II

```
1  module type VALUE = sig          13  val constant : constant -> typ -> t
2  type t                           14  val binop : operator -> typ -> t -> typ ->
3                                    15    t -> typ -> t
4  val name : string                16  val filter : bool -> typ -> t -> t
5                                    17
6  val bottom: t                    18  val backward_binop : operator -> typ ->
7  val top: t                       19    t -> typ -> t -> typ -> t -> t * t
8                                    20  val compare : operator -> bool -> typ ->
9  val subset: t -> t -> bool       21    t -> typ -> t -> (t * t)
10 val join: t -> t -> t            22
11 val meet: t -> t -> t            23  val print: printer -> t -> unit
12 val widen: 'a ctx -> t -> t -> t 24  end
```

Implementations for intervals, congruences, powerset of integers

They all abstract the same object

17

# Relational domains

```
1  // Hyp: a array, of size l ∈ [10, 20]
2  s = 0;
3  for(int i = 0; i < l; i++) {
4    s += a[i];        i ∈ [0, 20], l ∈ [10, 20], unable to prove safe access ✗
5  }
```

## The cartesian abstraction breaks relationality

$$\mathcal{P}(\mathcal{V} \to \mathbb{Z}) \xleftrightarrow[\alpha]{\gamma} \mathcal{V} \to \mathcal{P}(\mathbb{Z})$$

$$(\gamma \circ \alpha)(\{\, 0 \le x < y \le 2 \,\}) = (\gamma \circ \alpha)(\{\, (0,1); (0,2); (1,2) \,\})$$

$$= \gamma\,(x \mapsto \{\, 0, 1 \,\}, y \mapsto \{\, 1, 2 \,\})$$

$$= \{\, (0,1); (0,2); (1,1); (1,2) \,\}$$

$$= \{\, 0 \le x < 2, 0 < y \le 2 \,\}$$

18

## Relational domains - II

Mopsa relies on the Apron library [JM09], providing among others:

▶ Polyhedra [CH78], variants with PPL[BHZ08] or PPLite[BZ20]  $\sum_i \alpha_i V_i \leq \beta_i$
▶ Octagons [Min06b]  $\pm V_i \pm V_j \leq c_{i,j}$
▶ Grids [Bag+06]  $\sum_i \alpha_i V_i \equiv \beta_i[n]$

<div style="text-align: center;">Motivational example, with polyhedra</div>

```
1  // Hyp: a array, of size l ∈ [10, 20]
2  s = 0;
3  for(int i = 0; i < l; i++) {
4    s += a[i];                          0 ≤ i < l, ✔
5  }
```

# Architecture of Mopsa

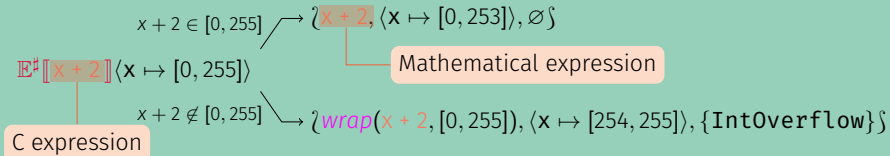Leveraging relational abstract domains

# Machine Numbers

▶ Numerical domains rely on **mathematical** numbers.

▶ C uses **finite precision** numbers with modular arithmetics.

▶ MACHINENUM lifts C statements to a math semantic.

## Dynamic Lifting

Consider a variable x declared as unsigned char.



$x + 2 \in [0, 255]$ ⟶ ⦉x + 2, ⟨x ↦ [0, 253]⟩, ∅⦊

Mathematical expression

$\mathbb{E}^\sharp[\![x + 2]\!]\langle x \mapsto [0, 255]\rangle$

C expression

$x + 2 \notin [0, 255]$ ⟶ ⦉*wrap*(x + 2, [0, 255]), ⟨x ↦ [254, 255]⟩, {IntOverflow}⦊

N.B: expression evaluation required here

# Machine Numbers – II

```
1  let eval exp man flow =
2    match exp.ekind with
3    | E_var v -> Some (Cases.singleton exp flow)          Variables do not overflow
4    | E_binop(op, e1, e2) ->
5      man.eval e1 flow >>$ fun n1 flow ->                  Evaluate e1 and bind each case ⟨n1⟩
6      man.eval e2 flow >>$ fun n2 flow ->                  Evaluate e2 and bind each case ⟨n2⟩
7      let vmin, vmax = rangeof exp.etyp in
8      let nexp = mk_binop n1 op n2 in
9      let ret = assume (mk_in nexp vmin vmax) man flow      Partition on condition
10         ~fthen:(fun flow ->                               nexp ∈ [vmin, vmax]
11             let flow = safe_c_integer_overflow flow in
12             Cases.singleton nexp flow)
13         ~felse:(fun flow ->
14             let nexp' = mk_wrap nexp vmin vmax in
15             let flow = raise_alarm IntegerOverflow flow in
16             Cases.singleton nexp' flow
17           ) in Some ret
18    | _ -> None
```

21

# Abstracting containers (strings, arrays) lengths

Consider a variable-length container $a$.

> Motivational example, with polyhedra

```
1  // Hyp: a container
2  s = 0;
3  for(int i = 0; i < container_length(a); i++) {
4    s += a[i];
5  }
```

We track its length through the introduction of a ghost numerical variable $\underline{len}(a)$

The relational domain will be able to infer relationships between $i$ and $\underline{len}(a)$.

N.B: In a non-relational setting, we could track values directly.

Convention: ghost variables are underscored.

# Abstracting containers (strings, arrays) lengths – II

Universal.Toy.String_length

```
1  let exec stmt man flow = let range = srange stmt in match skind stmt with
2    | S_assign ({ekind = E_var (s, _); etyp=T_string}, e) ->
3      Some (man.exec (mk_assign (mk_len_string_var s range)
4                      (mk_expr (E_len e) range) range) flow)
5
6    | S_assign ({ekind = E_subscript ({ekind = E_var (s, _)}, i); e) ->
7      Some (
8        assume (mk_in i (mk_zero range) (mk_len_string_var s range) range)
9                man flow
10               ~fthen:(safe_subscript_access_check)
11               ~felse:(fun flow ->
12                   let flow = raise_alarm invalid subscript access flow in
13                   Cases.empty flow))
14
15   | _ -> None
```

Case s = e — Rewrite and delegate into $\underline{len}(s) = len(e)$

Case s[i] = e

Case disjunction on $0 \leq i \leq \underline{len}(s)$

Safe, nothing to do

Return ⊥, with alarm metadata

23

## Abstracting containers (strings, arrays) lengths – III

```
1 string s;
2 if (rand(0, 1)) { s = "abcd"; }
3 else { s = "ab"; }
```

▶ Intervals $\qquad \underline{\text{len}}(s) \in [2, 4]$

▶ Intervals $\wedge$ Congruences $\quad \underline{\text{len}}(s) \in [2, 4] \wedge 2\mathbb{Z}$

▶ Powerset of integers $\qquad \underline{\text{len}}(s) \in \{2, 4\}$

```
1 string s = rand();
2 string t = s + s;
```

▶ Intervals $\underline{\text{len}}(s) \in [0, +\infty], \underline{\text{len}}(t) \in [0, +\infty]$

▶ Polyhedra $\qquad 0 \leq t = 2 \cdot \underline{\text{len}}(s)$

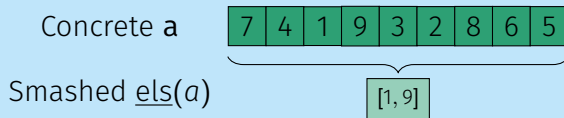NB: in case of dynamic allocations, ghost variables are a bit more complicated

# Abstracting containers (strings, arrays) contents

## Variable-length containers have unbounded size

Not abstracting their contents $\implies$ Non-terminating analysis

## The "smashing abstraction"

Idea: summarize every concrete cell of the container into an abstract one

Concrete $a$ | 7 | 4 | 1 | 9 | 3 | 2 | 8 | 6 | 5 |

Smashed $\underline{els}(a)$ | $[1, 9]$ |

## Weak update

Analyzing $a[2] = 12$, assuming $\underline{els}(a) \in [1, 9]$

✗ $\underline{els}(a) = [12, 12]$, as $\underline{els}(a)$ represents multiple concrete elements!

✓ $\underline{els}(a) \overset{\text{weak}}{=} [12, 12] \overset{\text{def}}{=} \underline{els}(a) \sqcup [12, 12] = [1, 12]$

$\mathbb{S}^{\sharp}[\![ c \overset{\text{weak}}{=} e ]\!] s \overset{\text{def}}{=} s \sqcup \mathbb{S}^{\sharp}[\![ c = e ]\!](s)$

# Abstracting containers (strings, arrays) contents – II

### What about weak read?

$r$ = $a$[3], assuming $\underline{els}(a) \in [1, 9]$.

$r \in [1, 9]$, and similarly with <u>non-relational domains</u>.
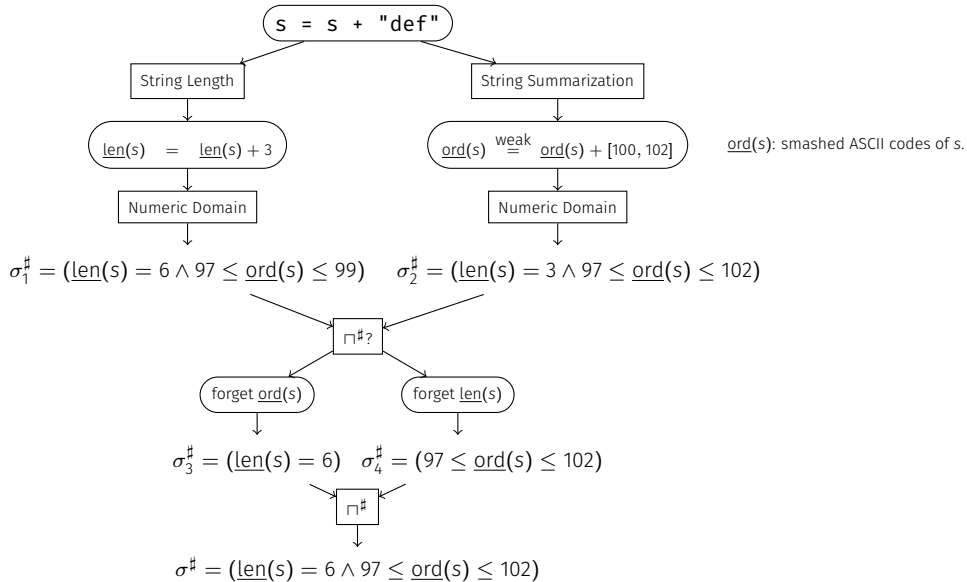
### Weak read and relational domains

$r$ = $a$[3], assuming $1 \leq \underline{els}(a) \leq i$.

✗ $1 \leq r = \underline{els}(a) \leq i$, as $\underline{els}(a)$ represents multiple concrete elements!

✓ $\mathbb{S}^{\sharp}[\![\, \text{expand}(\underline{els}(a), r) \,]\!](1 \leq \underline{els}(a) \leq i) = 1 \leq \underline{els}(a) \leq i \wedge 1 \leq r \leq i$

Intuitively: expand copies constraints. Cf. [Gop+04]

# The perils of reduced products with shared relational domains



$$s = s + \texttt{"def"}$$

String Length

$$\underline{len}(s) = \underline{len}(s) + 3$$

Numeric Domain

String Summarization

$$\underline{ord}(s) \stackrel{weak}{=} \underline{ord}(s) + [100, 102]$$

$\underline{ord}(s)$: smashed ASCII codes of $s$.

Numeric Domain

$$\sigma_1^\sharp = (\underline{len}(s) = 6 \wedge 97 \leq \underline{ord}(s) \leq 99) \quad \sigma_2^\sharp = (\underline{len}(s) = 3 \wedge 97 \leq \underline{ord}(s) \leq 102)$$

$\sqcap^\sharp$?

forget $\underline{ord}(s)$

forget $\underline{len}(s)$

$$\sigma_3^\sharp = (\underline{len}(s) = 6) \quad \sigma_4^\sharp = (97 \leq \underline{ord}(s) \leq 102)$$

$\sqcap^\sharp$

$$\sigma^\sharp = (\underline{len}(s) = 6 \wedge 97 \leq \underline{ord}(s) \leq 102)$$

Mopsa relies on <u>rewriting</u>, <u>symbolic expressions</u> and <u>ghost variables</u>

to leverage relational domains.

**The great power of relational domains comes with**

- ▶ Computational cost more than $\mathcal{O}(|\mathcal{V}|^3)$
- ▶ Force cohabitation of variables (cornerstone of Mopsa's design)
- ▶ Weak variables need specific operators (`expand`)
- ▶ Reduced product with sharing needs `merge` operator

# Architecture of Mopsa

Domains & their combinators

# Hasse diagram of domains



29

# Combinators

## Stacked Product

- ▶ Statements/expressions are dispatched to both sides
- ▶ Effects are collected to merge results soundly
- ▶ Reductions can be defined after evaluations of <u>expressions</u> or <u>statements</u>.

## Switch

- ▶ Pointwise lattice lifting
- ▶ For $\mathtt{tf} \in \{\,\mathbf{eval}, \mathbf{exec}\,\}$:

```
1  Switch(D1, D2).tf args =
2  let o_r = D1.tf args in
3  match o_r with
4   None -> D2.tf args
5   Some r -> r
```

## Compose

- ▶ Similar to switch
- ▶ Lattice operators can trigger operations on underlying domains

# Architecture of Mopsa

Transparency in static analysis

```
$ static-analysis-tool file
...
No errors found
```

What has been checked? What has not?

# Mopsa's approach to being transparent – implementation

C.Memory.Machine_Numbers

```
1  let eval exp man flow =
2    match exp.ekind with
3    | E_var v -> Some (Cases.singleton exp flow)
4    | E_binop(op, e1, e2) ->
5      man.eval e1 flow >>$ fun n1 flow ->
6      man.eval e2 flow >>$ fun n2 flow ->
7      let vmin, vmax = rangeof exp.etyp in
8      let nexp = mk_binop n1 op n2 in
9      let ret = assume (mk_in nexp vmin vmax) man flow
10         ~fthen:(fun flow ->
11             let flow = safe_c_integer_overflow flow in
12             Cases.singleton nexp flow)
13         ~felse:(fun flow ->
14             let nexp' = mk_wrap nexp vmin vmax in
15             let flow = raise_alarm IntegerOverflow flow in
16             Cases.singleton nexp' flow
17           ) in Some ret
18    | _ -> None
```

Transparent: mark that
e1 + e2 does not overflow

Standard: report alarm
Stored as metadata in 'a flow

## Mopsa's approach to being transparent

- ▶ Reporting status of all proofs / checks in every analyzed context
- ▶ Quantitative precision measure

$$\text{Selectivity} = \frac{\#\text{checks proved safe}}{\#\text{checks}}$$

```
1  int main() {
2    int n;
3    int y = -1;
4    for(int x = 0; x < n; x++)
5      y++;
6  }
```

| Stmt | Itv | Poly |
|------|-----|------|
| x++ | Safe | Safe |
| y++ | Alarm | Safe |
| Selectivity | 50% | 100% |

33

## Benefits of the approach

▶ Easy to implement

▶ "2,756 alarms" ⤳ 87% checks proved correct – "selectivity"

▶ ~~Program size~~ ⤳ "expression complexity"

Analysis of coreutils `fmt`

```
Checks summary: 21247 total, ✔ 18491 safe, ✗ 129 errors, △ 2627 warnings
  Stub condition: 690 total, ✔ 513 safe, ✗ 3 errors, △ 174 warnings
  Invalid memory access: 8139 total, ✔ 7142 safe, ✗ 4 errors, △ 993 warnings
  Division by zero: 499 total, ✔ 445 safe, △ 54 warnings
  Integer overflow: 11581 total, ✔ 10177 safe, △ 1404 warnings
  Invalid shift: 163 total, ✔ 163 safe
  Invalid pointer comparison: 37 total, ✗ 37 errors
  Invalid pointer subtraction: 85 total, ✗ 85 errors
  Insufficient variadic arguments: 1 total, ✔ 1 safe
  Insufficient format arguments: 26 total, ✔ 25 safe, △ 1 warning
  Invalid type of format argument: 26 total, ✔ 25 safe, △ 1 warning
```

34

# Mopsa's approach to being transparent – soundness assumptions

## Soundness assumptions, through an example

`extern f(int *x)`, handling gradations

1. Crash ✗
2. Ignore silently ✗
3. Assume and report: f has no effect
4. Assume and report: f has any effect on its parameters
5. Assume and report: f has any effect on its parameters and on globals

Related: soundiness paper [Liv+15]

# Current analyses in Mopsa

Control-flow tokens

## Control-flow tokens

- ▶ Astrée [Cou+06, footnote 4, page 6], Mopsa do not iterate over a CFG
- ▶ Proceeds by induction over the syntax
- ▶ Non-local control-flow is represented by <u>control-flow tokens</u>
  *cur* is the standard control-flow, *brk* states interrupted by a **break**
- ▶ Also applies to **goto**, **return**, **raise**, …
- ▶ Implementation: `'a flow` contains a `'a TokenMap.t`

```
1  int i = 0;
2  while(i < 100) {
3    i = i + rand(1, 3);
4    if(i == 42) break;
5  }
6
```

$cur \mapsto i \in [1, 99]$
$cur \mapsto i \in [2, 102]$
$brk \mapsto i \in [42, 42]$

$cur \mapsto i \in ([2, 102] \sqcap [100, +\infty]) \sqcup [42, 42]$

More details: [Mon21, section 2.4.3]

# Control-flow tokens – II

```
1  let exec stmt man flow = match stmt.skind with
2    | S_break ->
3      Some (Cases.return () (Flow.rename T_cur T_break man.lattice flow))
4
5    | S_while (cond, body) ->
6      Some (
7        lfp man cond body (Flow.rm T_break flow) flow >>% fun lfp_flow ->
8        man.exec (mk_assume (mk_not cond)) lfp_flow >>% fun lfp_flow ->
9        let lfp_flow = Flow.add T_cur (Flow.get T_break man.lattice lfp_flow)
10                       man.lattice lfp_flow in
11       let lfp_flow = Flow.set T_break (Flow.get T_break man.lattice flow)
12                       man.lattice lfp_flow in
13       Cases.return () lfp_return
14     )
15
16   | _ -> None
```

Standard flow transferred to *brk*

At loop end:
join *brk* and *cur* into *cur*

37

# Current analyses in Mopsa

## C analysis

▶ Checks for run-time errors
  (integer overflows, invalid dereferences, …)

▶ Supports ints, floats, pointers, structs, …

▶ Inlining-based analysis
  ⚠ scalability

▶ No concurrency support

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

1 Each pointer is mapped to the set of pointed bases
(e.g, variables or dynamically allocated memory)

2 Offsets are ghost numerical variables: offset(p)
Can express relations between offsets and numeric variables:

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

# Cells domain

A low-level memory abstraction [Min06a]

- ▶ Translates memory accesses into scalar accesses
- ▶ `type cell = { base: var; offset: Z.t; typ: scalar_type }`
- ▶ Synthesis for memory dereferences
- ▶ Supports type-punning, pointer arithmetic. No bitfield support implemented

Lipari-themed example (little-endian case)

```
1 uint32_t eax = 0xF0CACC1A;
2 uint8_t x = *((uint8_t *) &eax + 3);
```

$\mathbb{S}^\sharp [\![ \text{uint32\_t eax = 0xF0CACC1A} ]\!]$

$\quad \mathbb{S}^\sharp [\![ \underline{\text{cell}}(\wr\text{eax}, 0, u32\wr) = \text{0xF0CACC1A} ]\!]$

$\mathbb{S}^\sharp [\![ \text{uint8\_t x = *((uint8\_t *) \&eax + 3);} ]\!]$

$\quad \mathbb{E}^\sharp_{cells} [\![ \text{*((uint8\_t *) \&eax + 3)} ]\!]$

$\qquad \mathbb{S}^\sharp [\![ \underline{\text{cell}}(\wr\text{eax}, 3, u8\wr) = \text{0xF0} ]\!]$

$\quad \leftarrow \underline{\text{cell}}(\wr\text{eax}, 3, u8\wr)$

$\quad \mathbb{S}^\sharp [\![ x = \underline{\text{cell}}(\wr\text{eax}, 3, u8\wr) ]\!]$

40

# String length domain [JMO18]

New ghost variables:

- <u>bytes</u>(`var`) size in bytes of memory block
- <u>slen</u>(`var`) position of first must 0.

```
1  void strcpy(char *dst, char *src) {
2    char *p = dst, *q = src;
3    while(*q != '\0') {
4      *p = *q; p++; q++;
5    }
6    *p = *q;
7  }
```

# String length domain

Conjunction of pre-conditions · The switch utility · Continuation (with pre-filtered state)

```
1  val switch : (expr list * ('a Flow.flow -> ('a,'r) cases)) list) ->
2      ('a, 'b) man -> 'a flow -> ('a, 'r) cases
```

Transfer function of `base[offset] = rhs`

```
1  switch [
2    (* set0 case *)
3    (* Offset condition: offset ∈ [0, length] *)
4    (* RHS condition: rhs = 0 *)
5    (* Transformation: length := offset; *)
6    [ mk_in offset zero length range;
7      mk_eq rhs zero range ],
8    man.exec (mk_assign length offset range)
9    ;
10
11   (* setn0 case *)
12   (* Offset condition: offset = length *)
13   (* RHS condition: rhs ≠ 0 *)
14   (* Transformation: length := [offset + 1, size]; *)
15   [ mk_eq offset length range;
16     mk_ne rhs zero range ],
17   assign_length_interval (add offset one range) size
18   ;
20   (* First unchanged case *)
21   (* Offset condition: offset ∈ [0, length - 1] *)
22   (* RHS condition: rhs ≠ 0 *)
23   (* Transformation: nop; *)
24   [ mk_in offset zero (pred length range) range;
25     mk_ne rhs zero range ],
26   (fun flow -> Post.return flow)
27   ;
28
29   (* Second unchanged case *)
30   (* Offset condition: offset >= length + 1 *)
31   (* RHS condition: ⊤ *)
32   (* Transformation: nop; *)
33   [ mk_ge offset (succ length range) range ],
34   (fun flow -> Post.return flow)
35
36 ]
37   man flow
```

42

# libc stubs

▶ Annotation of Libc through a contract language [OM20]

▶ Inspired from Frama-C ACSL

▶ Contract language not restricted to C

strlen contract

```
1  /*$
2   * requires: valid_string_or_fail(__s);
3   * ensures : return ∈ [0, size(__s) - 1];
4   * ensures : __s[return] == 0;
5   * ensures : ∀ int i ∈ [0, return - 1]: __s[i] != 0;
6   */
7  size_t strlen(const char *__s);
```

User-defined predicate, including
$\exists \text{int } i \in [0, \text{size}(\_\_s)-1] : \_\_s[i] == 0$

Quantifier interpretation: delegated to domains

# Some benchmarks

See SV-Comp 2024 results.

| Benchmark | # Tests | Total LOC | Time | Precision |
|-----------|---------|-----------|------|-----------|
| CWE121 | 2,508 | 234,930 | 3,064s | 22.13% |
| CWE122 | 1,556 | 166,664 | 1,948s | 25.84% |
| CWE124 | 758 | 93,372 | 961s | 36.94% |
| CWE126 | 600 | 75,984 | 769s | 46.83% |
| CWE127 | 758 | 89,022 | 963s | 37.07% |
| CWE190 | 3,420 | 440,749 | 4,356s | 78.13% |
| CWE191 | 2,622 | 340,884 | 3,236s | 78.87% |
| CWE369 | 497 | 83,238 | 674s | 70.42% |
| CWE415 | 190 | 17,990 | 228s | 100.00% |
| CWE416 | 118 | 14,782 | 142s | 67.80% |
| CWE469 | 18 | 1,520 | 22s | 100.00% |
| CWE476 | 216 | 20,427 | 254s | 100.00% |

Table 1: Juliet benchmarks (non-relational configuration, no partitioning).

| Benchmark | Time | Selectivity | # checks |
|-----------|------|-------------|----------|
| basename | 33.79s | 98.65% | 11,731 |
| comm | 42.67s | 97.32% | 12,654 |
| dircolors | 34.82s | 99.74% | 20,062 |
| dirname | 21.68s | 99.61% | 11,307 |
| echo | 19.26s | 99.43% | 11,010 |
| false | 14.50s | 99.72% | 10,774 |
| getlimits | 34.62s | 98.54% | 11,711 |
| hostid | 18.05s | 99.65% | 11,303 |
| id | 32.69s | 99.04% | 12,338 |
| link | 23.03s | 99.52% | 11,572 |
| logname | 20.36s | 99.66% | 11,307 |
| mkfifo | 34.87s | 99.20% | 11,807 |

Table 2: `coreutils` benchmarks (fully symbolic arguments, relational analysis).

# Current analyses in Mopsa

Python analysis

- ▶ Detects uncaught exceptions
- ▶ Type and value analyses available
- ▶ Supports crazy Python dynamic typing, semantics, …
- ▶ Does not support GC finalizers, `async`, `eval`

## Nominal types: classes, MRO

Pointer-like domain

$\mathcal{V} \to \mathcal{P}(\mathsf{Addr}^\sharp \cup \{\, \mathtt{LocUndef}, \mathtt{GlobUndef} \,\})$

## Structural types: attributes

Attribute abstraction + ghost variables

$\mathsf{Addr}^\sharp \to \mathsf{ObjS}^\sharp$

Fspath (from standard library)

```
1  class Path:
2    def __fspath__(self): return 42
3
4  def fspath(p):
5    if isinstance(p, (str, bytes)):
6      return p
7    elif hasattr(p, "__fspath__"):
8      r = p.__fspath__()
9      if isinstance(r, (str, bytes)):
10       return r
11   raise TypeError
12
13 fspath("/dev" if random() else Path())
```

# Attribute abstraction

## Using an under and an over-approximation

$\mathsf{ObjS}^{\sharp} = \{\, (l, u) \mid l \in \mathcal{P}(string), u \in \mathcal{P}(string) \cup \{\, \top \,\}, l \subseteq u \vee u = \top \,\}$

## Concretization

$$\gamma^{\sharp}_{\mathsf{ObjS}} : \begin{cases} \mathsf{ObjS}^{\sharp} & \to & \mathcal{P}(\mathcal{P}(string)) \\ (l, \top) & \mapsto & \{\, s \in \mathcal{P}(string) \mid l \subseteq s \,\} \\ (l, u) & \mapsto & \{\, s \in \mathcal{P}(string) \mid l \subseteq s \subseteq u \,\} \end{cases}$$

## Example

$\gamma^{\sharp}_{\mathsf{ObjS}}(\{\, a \,\}, \{\, a, b, c \,\}) = \{\, \{\, a \,\}, \{\, a, b \,\}, \{\, a, c \,\}, \{\, a, b, c \,\} \,\}$

## The recency abstraction [BR06]

▶ Precise analysis of object initialization

▶ Twofold partitioning:
  • by allocation site $l \in$ **Loc**
  • through a recency criterion: $(l, \boldsymbol{r})$ most recent allocation (with strong updates)
    $\phantom{through a recency criterion:}$ $(l, \boldsymbol{o})$ older addresses (summarized)

▶ Initially designed for analysis of low-level code (binaries, C)

▶ Also used in Type Analysis for JavaScript [JMT09]

```python
class Task:
    def __init__(self, weight):
        if weight < 0: raise ValueError
        self.weight = weight

l = [Task(2), Task(1), Task(4), Task(5)]
```

## Return of ghost variables

Composed on top of address, for attribute "weight":

$@^\sharp(\text{Task}, r) \cdot \text{weight} \mapsto [2, 2]$

$$@^\sharp(\text{Task}, r) \cdot \text{weight} \mapsto [2, 2]$$

$$\begin{cases} @^\sharp(\text{Task}, r) \cdot \text{weight} \mapsto [1, 1] \\ @^\sharp(\text{Task}, o) \cdot \text{weight} \mapsto [2, 2] \end{cases}$$

$$\begin{cases} @^\sharp(\text{Task}, r) \cdot \text{weight} \mapsto [4, 4] \\ @^\sharp(\text{Task}, o) \cdot \text{weight} \mapsto [1, 2] \end{cases}$$

$$\begin{cases} @^\sharp(\textbf{Task}, \textbf{\textit{r}}) \cdot \textbf{weight} \mapsto [\textbf{5}, \textbf{5}] \\ @^\sharp(\text{Task}, o) \cdot \text{weight} \mapsto [1, 4] \end{cases}$$

# Recency abstraction – III

### Task creation

```python
1  class Task:
2    def __init__(self, weight):
3      if weight < 0: raise ValueError
4      self.weight = weight
5
6  m = [1, 2]
7  l = [Task(i) for i in m]
8  l.append(Task(3))
```

## Type analysis
Nominal types used in abstract addresses. No need for allocation-site in `Tasks`. But helpful for lists!

## Value analysis
Use allocation sites for `range` objects.

## Variable allocation policies

▶ Type-based (nominal) and/or location-based partitioning.

▶ Different configurations depending on type/value analysis.

| Name | LOC | Type Analysis | | Exceptions detected | | | Non-relational Value Analysis | | Exceptions detected | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | Time | Mem. | Type | Index | Key | Time | Mem. | Type | Index | Key |
| 🐍 nbody.py | 157 | | | | | | | | 0 | 1 | 1 |
| 🐍 scimark.py | 416 | | | | | | | | 1 | 0 | 0 |
| 🐍 richards.py | 426 | | | | | | | | 1 | 2 | 0 |
| 🐍 unpack_seq.py | 458 | | | | | | | | 0 | 0 | 0 |
| 🐍 go.py | 461 | | | | | | | | 3 | 20 | 0 |
| 🐍 hexiom.py | 674 | 1 | | | | | | | 0 | 21 | 3 |
| 🐍 regex_v8.py | 1792 | | | | | | | | 0 | 145 | 0 |
| 📘 processInput.py | 1417 | 1 | | | | | | | 7 | 4 | 1 |
| 📘 choose.py | 2562 | 1. | | | 22 | 7 | 2.9m | 3.7GB | 12 | 13 | 7 |
| Total | 9294 | 4.0m | 2.8GB | 59 | 2214 | 12 | 13m | 9.1GB | 59 | 228 | 12 |

**Conclusion**

The non-relational value analysis

- ▶ does not remove false type alarms
- ▶ significantly reduces index errors
- ▶ is ≃ 3× costlier

**Heuristic packing and relational analyses**

- ▶ Static packing, using function's scope
- ▶ Rules out all 145 alarms of 🐍 `regex_v8.py` (1792 LOC) at 2.5× cost

51

# Selectivity of the non-relational value analysis

| Name | Attributes | Types | Indexes | Keys | Values | Overflows | Divisions |
|---|---|---|---|---|---|---|---|
| 🐍 scimark.py | 746/746 | 844/844 | 2/5 | | 29/30 | 21/43 | 20/21 |
| 🐍 richards.py | 352/353 | 389/389 | 2/4 | | 2/3 | | 2/2 |
| 🐍 unpack_seq.py | 807/807 | 1210/1210 | | | 1/1 | | |
| 🐍 go.py | 664/697 | 728/728 | 2/20 | | 7/7 | 6/12 | 4/6 |
| 🐍 hexiom.py | 598/598 | 672/672 | 10/32 | 0/3 | 23/24 | | |
| 🐍 regex_v8.py | 7357/7357 | 8349/8349 | 1913/2057 | | 63/63 | | |
| 🔵 processInput.py | 617/619 | 790/792 | 12/12 | 0/1 | 0/1 | 2/2 | |
| 🔵 choose.py | 2519/2521 | 2997/2999 | 28/39 | 4/8 | 9/24 | 7/17 | |

Selectivity of the analysis on some classes of exceptions

Selectivity = Number of proved safe operations / Total number of checks

An empty cell denotes a program where the kind of exception cannot happen

# Current analyses in Mopsa

Python+C analysis

## Python+C analysis overview

### Assessment 20% of the 200 most popular Python libraries rely on C code

- ▶ Performance (numpy)
- ▶ System libraries (pygit2)

### Dangers

- ▶ Different values ($\mathbb{Z}$ vs. `Int32`)
- ▶ Shared memory state

### Our approach

- ▶ **Combined** analysis of C, Python and interface code
- ▶ Previous works [TM07; FF08; LLR20] : JNI ⚠ ⇝ Java, low precision

# Multilanguage code – example

### counter.c

```c
1   typedef struct {
2       PyObject_HEAD;
3       int count;
4   } Counter;
5
6   static PyObject*
7   CounterIncr(Counter *self, PyObject *args)
8   {
9       int i = 1;
10      if(!PyArg_ParseTuple(args, "|i", &i))
11          return NULL;
12
13      self->count += i;
14      Py_RETURN_NONE;
15  }
16
17  static PyObject*
18  CounterGet(Counter *self)
19  {
20      return Py_BuildValue("i", self->count);
21  }
```

### count.py

```python
22  from counter import Counter
23  from random import randrange
24
25  c = Counter()
26  power = randrange(128)
    c.incr(2**power-1)
28  c.incr()
    r = c.get()
```

▶ $\texttt{power} \leq 30 \Rightarrow \texttt{r} = 2^{\texttt{power}}$

▶ $\texttt{power} = 31 \Rightarrow r = -2^{31}$

▶ $32 \leq \texttt{power} \leq 64$: OverflowError:
   signed integer is greater than maximum

▶ $\texttt{power} \geq 64$: OverflowError:
   Python int too large to convert to C long

54

# High-level idea

## Difficulty: shared memory

- ▶ Each language may change the memory state, and has a different view of it
- ▶ Synchronization? We could perform a full state translation, but
  - the cost would be high in the analysis
  - some abstractions can be shared between Python and C

## State separation ⤳ reduced synchronization

- ▶ Observation: structures are directly dereferenceable by one language only
- ▶ Switch to other language otherwise (`c.incr()` ⤳ `self->count += 1`)
  Additional hypothesis: C accesses to Python objects through the API
- ▶ Synchronization: only when objects change language for the first time
- ▶ Mopsa supports shared abstractions

From distinct Python and C analyses. . . to a multilanguage analysis!

56

# Multilanguage analysis benchmarks

## Corpus selection

▶ Popular, real-world libraries available on GitHub, averaging 412 stars.

▶ Whole-program analysis: we use the tests provided by the libraries.

| Library | C + Py. Loc | Tests | $\bullet$/test | $\frac{\text{\# proved checks}}{\text{\# checks}}$ % | # checks |
|---|---|---|---|---|---|
| noise | 1397 | 15/15 | 1.2s | 99.7% | 6690 |
| cdistance | 2345 | 28/28 | 4.1s | 98.0% | 13716 |
| llist | 4515 | 167/194 | 1.5s | 98.8% | 36255 |
| ahocorasick | 4877 | 46/92 | 1.2s | 96.7% | 6722 |
| levenshtein | 5798 | 17/17 | 5.3s | 84.6% | 4825 |
| bitarray | 5841 | 159/216 | 1.6s | 94.9% | 25566 |

# Easing development

Developing sound, precise, scalable, static analyzers is a challenge!

▶ Investing time in creating specific tools can help a lot!
Examples in the remainder of this section.

▶ Reproducible science
Experimental results should be reproducible, through artifacts

▶ Career
  • tool dev. vs papers
  • software life and changing jobs

# Easing development

CI, tests & benchmarks

# CI, tests & benchmarks

## Detecting breaking changes using continuous integration

► `mopsa-diff` to compare with previous results

► Reusing all benchmarks from our experimental evaluations

## Benchmark selection

Our benchmarks are

► third-party real code
► open-source – for the sake of reproducible science
► unmodified*
  • Underscores practicality of our approach
  • * stubs can be added in marginal cases

# Comparing analysis reports

## `mopsa-diff` script

▶ compares analysis report(s): either single output or set of outputs
▶ usecases: different configurations, different versions of Mopsa

```
--- baseline/touch-many-symbolic-args-a4.json
+++ pplite/touch-many-symbolic-args-a4.json

- time: 589.0760
+ time: 675.1761

+ parse-datetime.y:1399.44-46: alarm: Invalid memory access
- parse-datetime.y:965.56-71: alarm: Invalid memory access
- parse-datetime.y:980.25-52: alarm: Invalid memory access
- parse-datetime.y:1003.23-50: alarm: Invalid memory access
- parse-datetime.y:921.56-71: alarm: Invalid memory access
- parse-datetime.c:1733.2-8: alarm: Invalid memory access
- parse-datetime.y:781.26-41: alarm: Invalid memory access
- parse-datetime.y:772.23-38: alarm: Invalid memory access
- parse-datetime.y:755.23-38: alarm: Invalid memory access
- parse-datetime.y:973.25-52: alarm: Invalid memory access
- parse-datetime.y:610.8-41: alarm: Invalid memory access
- parse-datetime.y:743.25-40: alarm: Invalid memory access
```

```
139 reports compared
avg. time change        +52.065s
avg. speedup               -36%
new alarms                    2
removed alarms               32
new assumptions               0
removed assumptions           0
new successes                 0
new failures                  0
```

# Easing development

Static analyzer interfaces

# Where static analyzers usually start from

▶ Analysis output                                                   Too coarse

▶ Printing abstract state using builtins             Not interactive

▶ Interpretation trace            Can be dozens of gigabytes of text

```
+ S [| set_program_name(argv[0]); |]
| | | + S [| add(argv0)
| | | |     argv0 = argv[0]; |]
| | | | + S [| add(argv0) |]
| | | | | + S [| add(argv0) |] in below(c.iterators.intraproc)
| | | | | | + S [| add(argv0) |] in C/Scalar
| | | | | | | + S [| add(offset{argv0}) |] in Universal
| | | | | | | o S [| add(offset{argv0}) |] in Universal done [0.0001s, 1 case]
| | | | | | o S [| add(argv0) |] in C/Scalar done [0.0001s, 1 case]
| | | | | | + S [| add(argv0) |] in below(c.memory.lowlevel.cells)
| | | | | | | + S [| add(offset{argv0}) |] in Universal
| | | | | | | o S [| add(offset{argv0}) |] in Universal done [0.0001s, 1 case]
| | | | | | o S [| add(argv0) |] in below(c.memory.lowlevel.cells) done [0.0001s, 1 case]
| | | | | o S [| add(argv0) |] in below(c.iterators.intraproc) done [0.0001s, 1 case]
| | | | o S [| add(argv0) |] done [0.0002s, 1 case]
| | | + S [| argv0 = argv[0]; |]
| | | | + S [| argv0 = (signed char *) @argv{0}:ptr; |] in below(c.iterators.intraproc)
| | | | | + S [| argv0 = (signed char *) @argv{0}:ptr; |] in C/Scalar
| | | | | | + S [| offset{argv0} = (offset{@argv{0}:ptr} + 0); |] in Universal
| | | | | | | + S [| offset{argv0} = (offset{@argv{0}:ptr} + 0); |] in below(universal.iterators.intraproc)
```

# An interactive engine acting as abstract debugger

GDB-like interface to the abstract interpretation of the program

## Demo!

▶ Breakpoints
- Program location
- Specific transfer function, analysis of subexpression
- Alarm: jumping <u>back</u> to statement generating first alarm

▶ Navigation

▶ Observation of the abstract state
- Full state
- Projection on specific variables

▶ Some scripting capabilities

# IDE support

- ▶ Language Server Protocol for linters (report alarms)
- ▶ Debug Adapter Protocol providing interactive engine interface
- ▶ Both protocols introduced by VSCode, supported by multiple IDEs

# Easing development

Plug-ins to observe the analysis

# Hooks: a plug-in system of analysis observers

## Hooks

Observe analyzer state      before/after any expression/statement analysis

## Current hooks

▶ Logs: trace of interpretation performed by the analysis

▶ Thresholds for widening

▶ Coverage

▶ Heuristic unsoundness/imprecision detection

▶ Profiling

## Coverage

► Global metric for the analysis' results

► Good to detect issues in the instrumentation of the fully context-sensitive analysis

### No symbolic argument

```
./src/coreutils-8.30/src/fmt.c:
    'main' 76% of 72 statements analyzed
    'set_prefix' 100% of 12 statements analyzed
    'same_para' 100% of 1 statement analyzed
    'get_line' 100% of 30 statements analyzed
    'fmt' 100% of 7 statements analyzed
    'base_cost' 100% of 16 statements analyzed
    'line_cost' 100% of 10 statements analyzed
    'get_prefix' 100% of 18 statements analyzed
```

### Symbolic arguments

```
./src/coreutils-8.30/src/fmt.c:
    'main' 100% of 72 statements analyzed
```

65

# Heuristic unsoundness/imprecision detection

### Detection of unsound transfer functions

Bottom shouldn't appear after some statements (such as assignments).

Goblint: ensure that at least one branch of a conditional is analyzed.

### Detection of imprecise analysis

Warns when assignments to top are performed

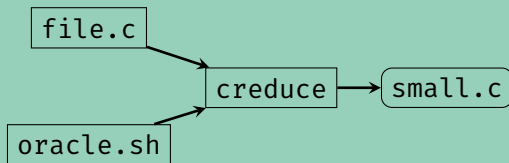Simplifies the search for sources of large imprecision

# Profiling

## Standard profiling

Measures which parts of Mopsa are the most time-consuming

## Abstract profiling hook

Measures which parts of the <u>analyzed program</u> are the most time-consuming

▶ Loop-level profiling

▶ Function-level profiling



Mopsa analysis of coreutils fmt

# Easing development

Testcase reduction

# Testcase reduction

## Motivation

▶ Static analyzers are complex piece of code and may contain bugs

▶ In practice, some bugs will only be detected in large codebases

▶ Debugging extremely difficult: size of the program, analysis time

## Automated testcase reduction using `creduce` [Reg+12]

## Testcase reduction – II

### Internal errors debugging

- ▶ Highly helpful to significantly reduce debugging time of runtime errors (Apron mishandlings, raised exceptions, …)
- ▶ Has been applied to coreutils programs, SV-Comp programs of 10,000+ LoC

### Differential-configuration debugging

```
$ mopsa-c -config=confA.json file.c
Alarm: assertion failure
$ mopsa-c -config=confB.json file.c
No alarm
```

Has been used to simplify cases in externally reported soundness issues

## Handling multi-file projects

### creduce limited to reducing a specific file

Mitigation: generate a pre-processed, standalone file

Painful operation on large projects such as `coreutils`

### Mopsa supports multi-file C projects

▶ `mopsa-build`
  - Records compiler/linker calls and their options
  - Creates a compilation database

  ⤳ `mopsa-build make` drop-in replacement for `make`

▶ `mopsa-c` leverages the compilation database

  ```
  mopsa-c mopsa.db -make-target=fmt
  ```

▶ Option to generate a single, preprocessed file

# Conclusion

## External fixpoint engine

- ▶ Mopsa: each iterator (loops, gotos, calls) defines its fixpoint computation.
- ▶ Alternative: unified, external fixpoint engine
- ▶ Used by Goblint team [Saa+24], Lermusiaux and Montagu [LM24].

## Systematic relationship between concrete and abstract domains

- ▶ See e.g, Michelland, Zakowski, and Gonnord [MZG24], Keidel and Erdweg [KE19]
- ▶ Lighter than a formally verified analyzer? [Jou+15]

Modular Open Platform for Static Analysis [Jou+19]
`gitlab.com/mopsa/mopsa-analyzer`

Goals: explore new designs, ease development of (relational) analyses

**One AST to rule them all**

- Multilanguage support
- Expressiveness
- Reusability

**Unified domain signature**

- Semantic rewriting
- Loose coupling
- Observability

**DAG of abstractions**

- Relational domains
- Composition
- Cooperation

Feedback wanted!

Anonymous survey, or come and talk to me!

# References – I

[Bag+06]    Roberto Bagnara et al. "Grids: A Domain for Analyzing the Distribution of Numerical Values". In: Lecture Notes in Computer Science. Springer, 2006, pp. 219–235.

[Bau+22]    Guillaume Bau et al. "Abstract interpretation of Michelson smart-contracts". In: ed. by Laure Gonnord and Laura Titolo. ACM, 2022, pp. 36–43. DOI: 10.1145/3520313.3534660.

[BHZ08]    Roberto Bagnara, Patricia M. Hill, and Enea Zaffanella. "The Parma Polyhedra Library: Toward a complete set of numerical abstractions for the analysis and verification of hardware and software systems". In: Sci. Comput. Program. 1-2 (2008), pp. 3–21.

## References – II

[BR06]     G. Balakrishnan and T. W. Reps. "Recency-Abstraction for Heap-Allocated Storage". In: LNCS. Springer, 2006, pp. 221–239.

[BZ20]     Anna Becchi and Enea Zaffanella. "PPLite: Zero-overhead encoding of NNC polyhedra". In: Inf. Comput. (2020), p. 104620. DOI: 10.1016/J.IC.2020.104620.

[CH78]     Patrick Cousot and Nicolas Halbwachs. "Automatic Discovery of Linear Restraints Among Variables of a Program". In: 1978.

[Cou+06]   Patrick Cousot et al. "Combination of Abstractions in the Astrée Static Analyzer". In: 2006.

## References – III

[DM19]     David Delmas and Antoine Miné. "Analysis of Software Patches Using Numerical Abstract Interpretation". In: ed. by Bor-Yuh Evan Chang. Lecture Notes in Computer Science. Springer, 2019, pp. 225–246. DOI: 10.1007/978-3-030-32304-2_12.

[DOM21]    David Delmas, Abdelraouf Ouadjaout, and Antoine Miné. "Static Analysis of Endian Portability by Abstract Interpretation". In: Lecture Notes in Computer Science. Springer, 2021, pp. 102–123.

[FF08]     Michael Furr and Jeffrey S. Foster. "Checking type safety of foreign function calls". In: ACM Trans. Program. Lang. Syst. (2008).

## References – IV

[Gop+04]   D. Gopan et al. "Numeric Domains with Summarized Dimensions". In: LNCS. Springer, 2004, pp. 512–529.

[JM09]   Bertrand Jeannet and Antoine Miné. "Apron: A Library of Numerical Abstract Domains for Static Analysis". In: Lecture Notes in Computer Science. Springer, 2009, pp. 661–667. DOI: 10.1007/978-3-642-02658-4_52.

[JMO18]   Matthieu Journault, Antoine Miné, and Abdelraouf Ouadjaout. "Modular Static Analysis of String Manipulations in C Programs". In: ed. by Andreas Podelski. Lecture Notes in Computer Science. Springer, 2018, pp. 243–262. DOI: 10.1007/978-3-319-99725-4_16.

[JMT09]    S. H. Jensen, A. Møller, and P. Thiemann. "Type Analysis for JavaScript". In: LNCS. Springer, 2009, pp. 238–255.

[Jou+15]   Jacques-Henri Jourdan et al. "A Formally-Verified C Static Analyzer". In: ed. by Sriram K. Rajamani and David Walker. ACM, 2015, pp. 247–259. DOI: 10.1145/2676726.2676966.

[Jou+19]   M. Journault et al. "Combinations of reusable abstract domains for a multilingual static analyzer". In: New York, USA, July 2019, pp. 1–17.

[KE19]     Sven Keidel and Sebastian Erdweg. "Sound and reusable components for abstract interpretation". In: Proc. ACM Program. Lang. OOPSLA (2019), 176:1–176:28. DOI: 10.1145/3360602.

[Liv+15]    Benjamin Livshits et al. "In defense of soundiness: a manifesto". In: Commun. ACM 2 (2015), pp. 44–46. DOI: 10.1145/2644805.

[LLR20]    Lee, Lee, and Ryu. "Broadening Horizons of Multilingual Static Analysis: Semantic Summary Extraction from C Code for JNI Program Analysis". In: 2020.

[LM24]    Pierre Lermusiaux and Benoît Montagu. "Detection of Uncaught Exceptions in Functional Programs by Abstract Interpretation". In: ed. by Stephanie Weirich. Lecture Notes in Computer Science. Springer, 2024, pp. 391–420. DOI: 10.1007/978-3-031-57267-8_15.

[MFM24]   Raphaël Monat, Aymeric Fromherz, and Denis Merigoux. "Formalizing Date Arithmetic and Statically Detecting Ambiguities for the Law". In: ed. by Stephanie Weirich. Lecture Notes in Computer Science. Springer, 2024, pp. 421–450. DOI: 10.1007/978-3-031-57267-8_16.

[Min06a]   Antoine Miné. "Field-sensitive value analysis of embedded C programs with union types and pointer arithmetics". In: ed. by Mary Jane Irwin and Koen De Bosschere. ACM, 2006, pp. 54–63. DOI: 10.1145/1134650.1134659.

[Min06b]   Antoine Miné. "The octagon abstract domain". In: High. Order Symb. Comput. (2006).

[MM24]     Marco Milanese and Antoine Miné. "Generation of Violation Witnesses by Under-Approximating Abstract Interpretation". In: ed. by Rayna Dimitrova, Ori Lahav, and Sebastian Wolff. Lecture Notes in Computer Science. Springer, 2024, pp. 50–73. DOI: 10.1007/978-3-031-50524-9_3.

[MOM20a]   R. Monat, A. Ouadjaout, and A. Miné. "Static Type Analysis by Abstract Interpretation of Python Programs". In: LIPIcs. 2020.

[MOM20b]   R. Monat, A. Ouadjaout, and A. Miné. "Value and allocation sensitivity in static Python analyses". In: ACM, 2020, pp. 8–13. DOI: 10.1145/3394451.3397205.

[MOM21]   R. Monat, A. Ouadjaout, and A. Miné. "A Multilanguage Static Analysis of Python Programs with Native C Extensions". In: 2021.

[Mon+24]   Raphaël Monat et al. "Mopsa-C: Improved Verification for C Programs, Simple Validation of Correctness Witnesses (Competition Contribution)". In: ed. by Bernd Finkbeiner and Laura Kovács. Lecture Notes in Computer Science. Springer, 2024, pp. 387–392. DOI: 10.1007/978-3-031-57256-2_26.

[Mon21]   Raphaël Monat. "Static Type and Value Analysis by Abstract Interpretation of Python Programs with Native C Libraries". PhD thesis. Sorbonne Université, France, 2021.

[MZG24]     Sébastien Michelland, Yannick Zakowski, and Laure Gonnord.
            "Abstract Interpreters: A Monadic Approach to Modular Verification".
            In: Proc. ACM Program. Lang. ICFP (Aug. 2024). DOI:
            10.1145/3674646.

[NP18]      Kedar S. Namjoshi and Zvonimir Pavlinovic. "The Impact of Program
            Transformations on Static Program Analysis". In: ed. by
            Andreas Podelski. Lecture Notes in Computer Science. Springer, 2018,
            pp. 306–325. DOI: 10.1007/978-3-319-99725-4_19.

[OM20]     A. Ouadjaout and A. Miné. "A Library Modeling Language for the Static Analysis of C Programs". In: ed. by David Pichardie and Mihaela Sighireanu. Lecture Notes in Computer Science. Springer, 2020, pp. 223–247. DOI: 10.1007/978-3-030-65474-0_11.

[PM24]     Francesco Parolini and Antoine Miné. "Sound Abstract Nonexploitability Analysis". In: ed. by Rayna Dimitrova, Ori Lahav, and Sebastian Wolff. Lecture Notes in Computer Science. Springer, 2024, pp. 314–337. DOI: 10.1007/978-3-031-50521-8_15.

[Reg+12]   John Regehr et al. "Test-case reduction for C compiler bugs". In: ed. by Jan Vitek, Haibo Lin, and Frank Tip. ACM, 2012, pp. 335–346. DOI: 10.1145/2254064.2254104.

[Saa+24]   Simmo Saan et al. "Goblint: Abstract Interpretation for Memory Safety and Termination - (Competition Contribution)". In: ed. by Bernd Finkbeiner and Laura Kovács. Lecture Notes in Computer Science. Springer, 2024, pp. 381–386. DOI: 10.1007/978-3-031-57256-2_25.

[TM07]   Gang Tan and Greg Morrisett. "Ilea: inter-language analysis across Java and C". In: 2007.

[VMM23]   Milla Valnet, Raphaël Monat, and Antoine Miné. "Analyse statique de valeurs par interprétation abstraite de programmes fonctionnels manipulant des types algébriques récursifs". In: ed. by Timothy Bourke and Delphine Demange. Praz-sur-Arly, France, Jan. 2023, pp. 211–242.