

MOPSA

Multilingual Static Analysis for Program Verification using Abstract Interpretation

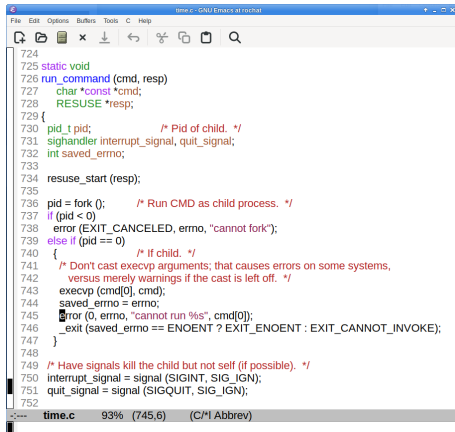
Guillaume Bau Jérôme Boillot David Delmas
Matthieu Journault Marco Milanese Antoine Miné
Raphaël Monat Abdelraouf Ouadjaout Francesco Parolini
Milla Valnet

APR team
LIP6, Sorbonne Université
Paris, France

SWHSec Conference
06/06/2024
Cyber Campus, Paris



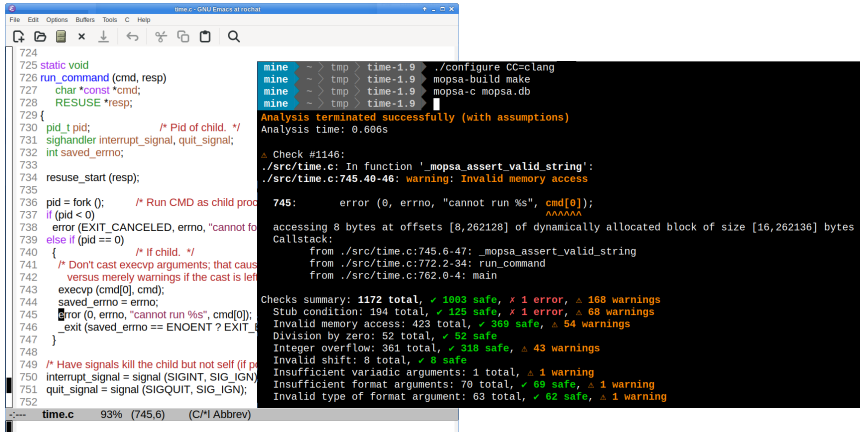
Program verification by static analysis



```
724
725 static void
726 run_command (cmd, resp)
727   char *const *cmd;
728   RESUSE *resp;
729 {
730   pid_t pid;          /* Pid of child. */
731   sighandler interrupt_signal, quit_signal;
732   int saved_errno;
733
734   resuse_start (resp);
735
736   pid = fork ();      /* Run CMD as child process. */
737   if (pid < 0)
738     error (EXIT_CANCELED, errno, "cannot fork");
739   else if (pid == 0)
740     { /* If child. */
741       /* Don't cast execvp arguments; that causes errors on some systems,
742          versus merely warnings if the cast is left off. */
743       execvp (cmd[0], cmd);
744       saved_errno = errno;
745       error (0, errno, "cannot run %s", cmd[0]);
746       _exit (saved_errno == ENOENT ? EXIT_ENOENT : EXIT_CANNOT_INVOKE);
747     }
748
749   /* Have signals kill the child but not self (if possible). */
750   interrupt_signal = signal (SIGINT, SIG_IGN);
751   quit_signal = signal (SIGQUIT, SIG_IGN);
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
```

- analyze **source code** at **compile time**, without executing it
- fully **automatic** (no annotation, no interaction), moderately efficient
- report on **non-functional correctness** and **assertion violations**
run-time errors, CWE, coding guidelines, etc.

Program verification by static analysis

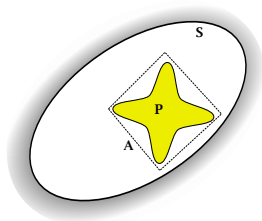


```
time.c - GNU Emacs at roshal
File Edit Options Buffers Tools C Help
724
725 static void
726 run_command(cmd, resp)
727 char *const *cmd;
728 RESUSE *resp;
729 {
730     pid_t pid;          /* Pid of child. */
731     sighandler_t interrupt_signal, quit_signal;
732     int saved_errno;
733
734     resuse_start(resp);
735
736     pid = fork();      /* Run CMD as child prog
737     if (pid < 0)
738         error (EXIT_CANCELED, errno, "cannot fo
739     else if (pid == 0)
740         /* if child. */
741         /* Don't cast execvp arguments; that caus
742         versus merely warnings if the cast is left
743         execvp (cmd[0], cmd);
744         saved_errno = errno;
745         error (0, errno, "cannot run %s", cmd[0]);
746         _exit (saved_errno == ENOENT ? EXIT_
747     }
748
749     /* Have signals kill the child but not self (if p
750     interrupt_signal = signal (SIGINT, SIG_IGN);
751     quit_signal = signal (SIGQUIT, SIG_IGN);
752
753     time.c 93% (745,6) (C/*I Abbrev)
```

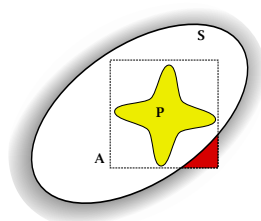
```
mine ~ > tmp > time-1.9 ./configure CC=clang
mine ~ > tmp > time-1.9 mopsa-build make
mine ~ > tmp > time-1.9 mopsa-c mopsa.db
mine ~ > tmp > time-1.9
Analysis terminated successfully (with assumptions)
Analysis time: 0.606s
△ Check #1146:
./src/time.c: In function '!_mopsa_assert_valid_string':
./src/time.c:745.40-46: warning: Invalid memory access
745:         error (0, errno, "cannot run %s", cmd[0]);
           ^^^^^^
accessing 8 bytes at offsets [8,262128] of dynamically allocated block of size [16,262136] bytes
Callstack:
  from ./src/time.c:745.6-47: !_mopsa_assert_valid_string
  from ./src/time.c:772.2-34: run_command
  from ./src/time.c:762.0-4: main
Checks summary: 1172 total, ✓ 1003 safe, ✗ 1 error, △ 168 warnings
Stub condition: 194 total, ✓ 125 safe, ✗ 1 error, △ 68 warnings
Invalid memory access: 423 total, ✓ 369 safe, △ 54 warnings
Division by zero: 52 total, ✓ 52 safe
Integer overflow: 361 total, ✓ 318 safe, △ 43 warnings
Invalid shift: 8 total, ✓ 8 safe
Insufficient variadic arguments: 1 total, △ 1 warning
Insufficient format arguments: 70 total, ✓ 69 safe, △ 1 warning
Invalid type of format argument: 63 total, ✓ 62 safe, △ 1 warning
```

- analyze **source code** at **compile time**, without executing it
- fully **automatic** (no annotation, no interaction), moderately efficient
- report on **non-functional correctness** and **assertion violations**
run-time errors, CWE, coding guidelines, etc.

Sound semantic static analysis



sound and precise analysis



false alarm

- compute an **over-approximation** A of **reachable program states** P and compare to a **specification** S
- $P \subseteq A$ by construction, guaranteed by abstract interpretation theory
- $A \subseteq S$: no bad state found \implies **guarantee the absence of error!**
- $A \not\subseteq S$: bad state found \implies **inconclusive** (true or false alarm)

A formal method **sound** and **incomplete** for the proof of absence of errors
 \neq SAST, linters, but also \neq model-checkers, deductive verification, theorem provers

Abstract Interpretation

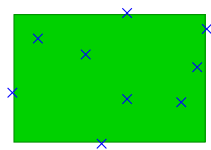
Representing reachable states

Abstract domains:

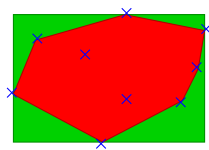
- computable, efficient **symbolic** representations of reachable states
⇒ non-representable sets are **over-approximated**



concrete memory states



box abstraction



polyhedral abstraction
relational abstraction

Abstract Interpretation

Efficiently computing over-approximations

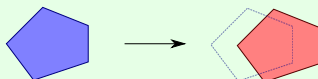
Abstract domains:

- computable, efficient **semantic operators** in the **abstract world** with soundness proof
- **iterators** traverse the program (CFG or AST)

Example: Polyhedra operations

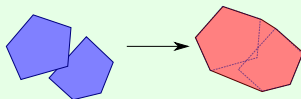
Assignments

• $X \leftarrow X + 1$ •
translation



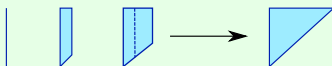
Branches: join

if ... then ... • else ... • fi •
convex hull



Loops: inductive invariants

while • ... do ... done
iteration with **widening** ▽



MOPSA

Goal: open-source static analysis platform

- for research and education in abstract interpretation
- easy to extend for new languages and new properties
- develop, test, reuse abstractions
- open-source (LGPL)

Contributors

Guillaume Bau Jérôme Boillot David Delmas
Matthieu Journault Marco Milanese Antoine Miné
Raphaël Monat Abdelraouf Ouadjaout
Francesco Parolini Milla Valnet

<https://gitlab.com/mopsa/mopsa-analyzer>

Some **classic** aspects:

- whole-program **forward abstract interpretation**
- by **induction** on the **AST**
- in a **collection** of **communicating abstract domains**, **sound** by design
- analyze **run-time errors** in **C** programs

More original aspects:

- new **languages**: **Python**, **OCaml**
- **multi-language**: programs **mixing C** and **Python**
- new **properties**: non-regression for **patches**, endianness **portability**, **exploitability**, inferring **counter examples**
- experiments in new analysis engineering and **architecture**...
- academic implementation! not industrial-scale! various levels of maturity!

Feature: Extensible syntax tree for multiple languages

Classic design: common intermediate representation (JVM, LLVM, Clight, ...)

- ✓ low-level: few and simple language constructs
- ✗ loss of high-level information and structure
- ✓ maximum sharing of abstract domains, easily retargetable (?)
- ✗ ... maybe not, if the language is too widely different! (Python?)

MOPSA design: extensible AST

- ✓ preserve the original AST of each source language
- ✓ and relevant fragments supported by well-known domains
e.g., scalar fragment, pointer-free fragment, integer arithmetic
- ✓ translation is dynamic, more power (exploit abstract information)
- ✗ less efficient
- ✗ hard to keep track of sub-languages, possible match failures at run-time

Feature: Extensible syntax tree for multiple languages

Classic design: common intermediate representation (JVM, LLVM, Clight, ...)

- ✓ low-level: few and simple language constructs
- ✗ loss of high-level information and structure
- ✓ maximum sharing of abstract domains, easily retargetable (?)
- ✗ ... maybe not, if the language is too widely different! (Python?)

MOPSA design: extensible AST

- ✓ preserve the original AST of each source language
- ✓ and relevant fragments supported by well-known domains
e.g., scalar fragment, pointer-free fragment, integer arithmetic
- ✓ translation is dynamic, more power (exploit abstract information)
- ✗ less efficient
- ✗ hard to keep track of sub-languages, possible match failures at run-time

Feature: Distributed iterators and delegation

```
type stmt_kind += S_while of expr * stmt
type stmt_kind += S_c_for of expr * expr * expr * stmt
type stmt_kind += S_py_for of expr * expr * stmt * stmt
```

simple loops: Universal.iterators.loops

- matches `S_while`
- computes a fixpoint with acceleration ∇

C loops: C.iterators.loops

- matches `S_c_for` (`init`, `cond`, `incr`, `body`)
- rewrite into `S_while`
- and call interpreter `recursively`

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

Python loops: Python.desugar.loops

- matches `S_py_for` (`target`, `iterable`, `body`)
- rewrite into `S_while`
- `optimize` simpler cases based on `dynamic types` (e.g., `ranges`)
- and call interpreter `recursively`

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

In MOPSA, iterators are just abstract domains... without internal abstract state

Feature: Distributed iterators and delegation

```
type stmt_kind += S_while of expr * stmt
type stmt_kind += S_c_for of expr * expr * expr * stmt
type stmt_kind += S_py_for of expr * expr * stmt * stmt
```

simple loops: `Universal.iterators.loops`

- matches `S_while`
- computes a fixpoint with acceleration ∇

C loops: `C.iterators.loops`

- matches `S_c_for` (`init`, `cond`, `incr`, `body`)
- rewrite into `S_while`
- and call interpreter `recursively`

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

Python loops: `Python.desugar.loops`

- matches `S_py_for` (`target`, `iterable`, `body`)
- rewrite into `S_while`
- `optimize` simpler cases based on `dynamic types` (e.g., `ranges`)
- and call interpreter `recursively`

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

In MOPSA, iterators are just abstract domains... without internal abstract state

Feature: Distributed iterators and delegation

```
type stmt_kind += S_while of expr * stmt
type stmt_kind += S_c_for of expr * expr * expr * stmt
type stmt_kind += S_py_for of expr * expr * stmt * stmt
```

simple loops: `Universal.iterators.loops`

- matches `S_while`
- computes a fixpoint with acceleration ∇

C loops: `C.iterators.loops`

- matches `S_c_for` (`init`, `cond`, `incr`, `body`)
- rewrite into `S_while`
- and call interpreter `recursively`

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

Python loops: `Python.desugar.loops`

- matches `S_py_for` (`target`, `iterable`, `body`)
- rewrite into `S_while`
- `optimize` simpler cases based on `dynamic types` (e.g., `ranges`)
- and call interpreter `recursively`

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

In MOPSA, iterators are just abstract domains... without internal abstract state

Feature: Dynamic expression rewriting with case analysis

Flexible, modular way for domains to communicate:

- **rewrite** (part of) expressions, introduce **new variables** (embed abstractions)
- **delegate** to other domains (unaware of other abstractions, decoupling)
- possible perform **case analysis** with disjunctions (split the abstract state)

Example: string length domain

Abstracts **character array** a as $P_a : a[v_a] = 0 \wedge (\forall k \in [0, v_a[: a[k] \neq 0)$

- maintain internally a map $a \mapsto P_a$
- introduce an integer variable v_a for each P_a
- rewrite a sub-expression $a[i] == 0$ in abstract element X^\sharp into a disjunction:
 - $(0, \mathbb{S}[i < v_a] X^\sharp)$
 - $(1, \mathbb{S}[i = v_a] X^\sharp)$ keep some relationality between i , $a[i]$, and v_a
 - $([0, 1], \mathbb{S}[i > v_a] X^\sharp)$
- $\mathbb{S}[i < v_a]$, etc. are handled by (relational) numeric abstract domains

Feature: Dynamic expression rewriting with case analysis

Flexible, modular way for domains to communicate:

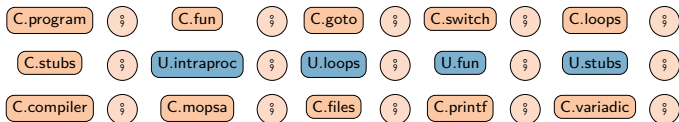
- **rewrite** (part of) expressions, introduce **new variables** (embed abstractions)
- **delegate** to other domains (unaware of other abstractions, decoupling)
- possible perform **case analysis** with disjunctions (split the abstract state)

Example: string length domain

Abstracts **character array** a as $P_a : a[v_a] = 0 \wedge (\forall k \in [0, v_a[: a[k] \neq 0)$

- maintain internally a map $a \mapsto P_a$
- introduce an integer variable v_a for each P_a
- rewrite a sub-expression $a[i] == 0$ in abstract element X^\sharp into a disjunction:
 - $(0, \mathbb{S}[i < v_a] X^\sharp)$
 - $(1, \mathbb{S}[i = v_a] X^\sharp)$ keep some relationality between i , $a[i]$, and v_a
 - $([0, 1], \mathbb{S}[i > v_a] X^\sharp)$
- $\mathbb{S}[i < v_a]$, etc. are handled by (relational) numeric abstract domains

Example: Domains for a C analysis



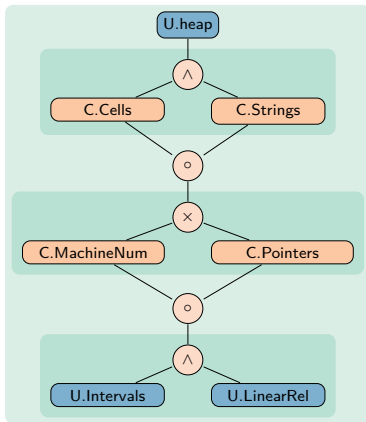
Sequence

Reduced product

Cartesian product

Universal

C specific



Requirement

For **soundness**, we need the C source **or a model** of **every function** called

stub for open

```
/*$
 * requires: exists int i in [0, size(__file) - 1]: __file[i] == 0;
 *
 * case "success":
 *   local: void* fd = new FileDescriptor;
 *   ensures: return == (int)fd;
 *
 * case "failure":
 *   assigns: _errno;
 *   ensures: return == -1;
 */
int open (const char *__file, int __oflag, ...);
```

- logic-based **contract language**, but with **C** expressions (sweet spot)
- stubs are **executed** at each call (not a specification to verify)
- for MOPSA: just another language (delegation, rewriting, etc.)
 - **add** abstract semantic operators for \forall , \exists , etc.
 - general notion of **resources** (generalizes dynamic memory allocation)

Requirement

For **soundness**, we need the C source **or a model** of **every function** called

stub for open

```
/*$  
 * requires: exists int i in [0, size(__file) - 1]: __file[i] == 0;  
 *  
 * case "success":  
 *   local: void* fd = new FileDescriptor;  
 *   ensures: return == (int)fd;  
 *  
 * case "failure":  
 *   assigns: _errno;  
 *   ensures: return == -1;  
 */  
int open (const char *__file, int __oflag, ...);
```

- logic-based **contract language**, but with **C** expressions (sweet spot)
- stubs are **executed** at each call (not a specification to verify)
- for MOPSA: just another language (delegation, rewriting, etc.)
 - **add** abstract semantic operators for \forall , \exists , etc.
 - general notion of **resources** (generalizes dynamic memory allocation)

C benchmarks: Juliet

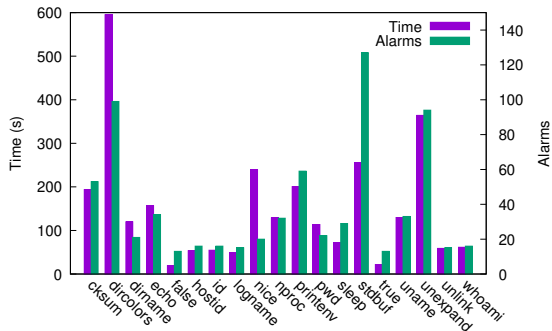
Experiments					
CWE	Lines	Time (h:m:s)	✓	⚠	
Stack-based Buffer Overflow	234k	00:59:12	89%	11%	
Heap-based Buffer Overflow	174k	00:37:12	86%	14%	
Buffer Underwrite	93k	00:18:28	86%	14%	
Buffer Over-read	75k	00:14:45	85%	15%	
Buffer Under-read	89k	00:18:26	87%	13%	
Integer Overflow	440k	01:24:47	52%	48%	
Integer Underflow	340k	01:02:27	55%	45%	
Divide By Zero	109k	00:13:17	55%	45%	
Double Free	17k	00:04:21	100%	0%	
Use After Free	14k	00:02:40	100%	0%	
Illegal Pointer Subtraction	1k	00:00:24	100%	0%	
NULL Pointer Dereference	21k	00:04:53	100%	0%	

Analyzed [12 CWEs](#) (13,261 tests) from [NIST Juliet v1.3](#)
each test comes with a good and a bad version

- ✓ good case safe **and** 1 error in bad case.
- ⚠ good case unsafe **or** many errors in bad case.
for all tests, the bad case reports an error (sound)

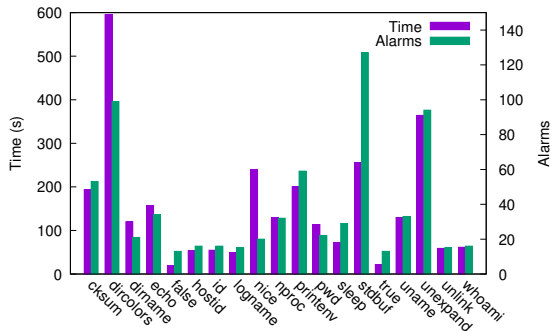
[with A. Ouadjaout @ SAS'20]

C benchmarks: Coreutils



- analyzed 19 programs from GNU Coreutils v8.30
- stub model for 1108 functions of the GNU libc
- `main` is called with a symbolic string array `argv` of arbitrary size

C benchmarks: Coreutils



- analyzed 19 programs from GNU Coreutils v8.30
- stub model for 1108 functions of the GNU libc
- `main` is called with a symbolic string array `argv` of arbitrary size

SV-COMP 2024 competition

1st place for Software Systems category !

Challenges in Python analysis

```
def f(a, b):  
    return a + b
```

```
 $\mathbb{E}[e_1 + e_2](f, \epsilon, \Sigma) = \text{def}$   
  let  $(f_1, \epsilon_1, \Sigma_1, v_1) = \mathbb{E}[e_1](f, \epsilon, \Sigma)$  in  
  let  $(f_2, \epsilon_2, \Sigma_2, v_2) = \mathbb{E}[e_2](f_1, \epsilon_1, \Sigma_1)$  in  
  if  $\text{hasattr}(v_1, \_radd\_\_)$ ,  $\Sigma_2$  then  
    let  $(f_3, \epsilon_3, \Sigma_3, v_3) = \mathbb{E}[v_1.\_radd\_(v_2)](f_2, \epsilon_2, \Sigma_2)$  in  
    if  $v_3 = \text{NotImpl} \wedge \text{typeof}(v_1) \neq \text{typeof}(v_2)$  then  
      if  $\text{hasattr}(v_2, \_radd\_\_)$ ,  $\Sigma_3$  then  
        let  $(f_4, \epsilon_4, \Sigma_4, v_4) = \mathbb{E}[v_2.\_radd\_(v_1)](f_3, \epsilon_3, \Sigma_3)$  in  
        if  $v_4 = \text{NotImpl}$  then  $\text{TypeError}(f_4, \epsilon_4, \Sigma_4)$  else  $(f_4, \epsilon_4, \Sigma_4, v_4)$   
      else  $\text{TypeError}(f_3, \epsilon_3, \Sigma_3)$   
    else if  $v_3 = \text{NotImpl}$  then  $\text{TypeError}(f_3, \epsilon_3, \Sigma_3)$  else  $(f_3, \epsilon_3, \Sigma_3, v_3)$   
  else if  $\text{hasattr}(v_2, \_radd\_\_)$ ,  $\Sigma_2$   $\wedge$   $\text{typeof}(v_1) \neq \text{typeof}(v_2)$  then  
    let  $(f_3, \epsilon_3, \Sigma_3, v_3) = \mathbb{E}[v_2.\_radd\_(v_1)](f_2, \epsilon_2, \Sigma_2)$  in  
    if  $v_3 = \text{NotImpl}$  then  $\text{TypeError}(f_3, \epsilon_3, \Sigma_3)$  else  $(f_3, \epsilon_3, \Sigma_3, v_3)$   
  else  $\text{TypeError}(f_2, \epsilon_2, \Sigma_2)$ 
```

- control-flow very dependent on the **dynamic** type
⇒ a precise **flow- and context-sensitive value** analysis is necessary!
- **complex** language to formalize!
using a functional big-step semantics, close to the abstract interpreter implementation
⇒ expression **rewriting** and **case analysis** is useful

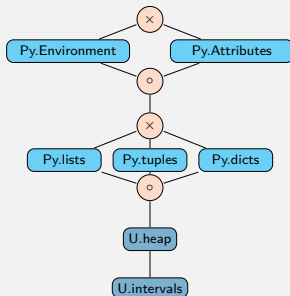
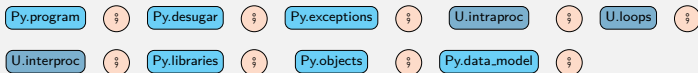
Challenges in Python analysis

```
def f(a, b):  
    return a + b
```

```
 $E[e_1 + e_2](f, \epsilon, \Sigma) = \text{def}$   
  let  $(f_1, \epsilon_1, \Sigma_1, v_1) = E[e_1](f, \epsilon, \Sigma)$  in  
  let  $(f_2, \epsilon_2, \Sigma_2, v_2) = E[e_2](f_1, \epsilon_1, \Sigma_1)$  in  
  if  $\text{hasattr}(v_1, \_radd\_, \Sigma_2)$  then  
    let  $(f_3, \epsilon_3, \Sigma_3, v_3) = E[v_1.\_radd\_(v_2)](f_2, \epsilon_2, \Sigma_2)$  in  
    if  $v_3 = \text{NotImpl} \wedge \text{typeof}(v_1) \neq \text{typeof}(v_2)$  then  
      if  $\text{hasattr}(v_2, \_radd\_, \Sigma_3)$  then  
        let  $(f_4, \epsilon_4, \Sigma_4, v_4) = E[v_2.\_radd\_(v_1)](f_3, \epsilon_3, \Sigma_3)$  in  
        if  $v_4 = \text{NotImpl}$  then  $\text{TypeError}(f_4, \epsilon_4, \Sigma_4)$  else  $(f_4, \epsilon_4, \Sigma_4, v_4)$   
      else  $\text{TypeError}(f_3, \epsilon_3, \Sigma_3)$   
    else if  $v_3 = \text{NotImpl}$  then  $\text{TypeError}(f_3, \epsilon_3, \Sigma_3)$  else  $(f_3, \epsilon_3, \Sigma_3, v_3)$   
  else if  $\text{hasattr}(v_2, \_radd\_, \Sigma_2) \wedge \text{typeof}(v_1) \neq \text{typeof}(v_2)$  then  
    let  $(f_3, \epsilon_3, \Sigma_3, v_3) = E[v_2.\_radd\_(v_1)](f_2, \epsilon_2, \Sigma_2)$  in  
    if  $v_3 = \text{NotImpl}$  then  $\text{TypeError}(f_3, \epsilon_3, \Sigma_3)$  else  $(f_3, \epsilon_3, \Sigma_3, v_3)$   
  else  $\text{TypeError}(f_2, \epsilon_2, \Sigma_2)$ 
```

- control-flow very dependent on the **dynamic** type
⇒ a precise **flow- and context-sensitive value** analysis is necessary!
- **complex** language to formalize!
using a functional big-step semantics, close to the abstract interpreter implementation
⇒ expression **rewriting** and **case analysis** is useful

Domains for Python value analysis



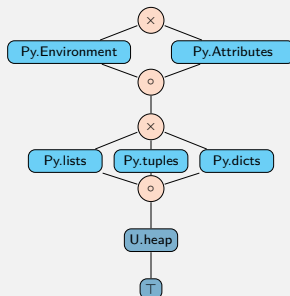
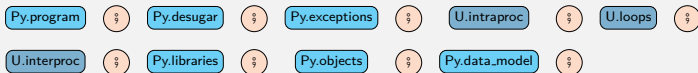
Example Python domain: Lists

- smash each list contents into a single “summary” variable
- keep the list length in a numeric variable

Application : analyze small programs (a few 100s of lines, few dependencies)

[with A. Fromherz & R. Monat @ SOAP'20, ECOOP'20]

Domains for Python type analysis



a type analysis is also easy to construct!

[with R. Monat @ SOAP'20, ECOOP'20]

Programs mixing C and Python

Python counter class in C

```
typedef struct {
    PyObject_HEAD;
    int counter;
} Counter;

static PyObject*
CounterIncr(Counter *self, PyObject *args) {
    int i = 1;
    if (!PyArg_ParseTuple(args, "|i", &i))
        return NULL;
    self->counter += i;
    Py_RETURN_NONE;
}

static PyObject* CounterGet(Counter *self) {
    return Py_BuildValue("i", self->counter);
}
```

Python client

```
from counter import Counter
from random import randrange

c = Counter()
power = randrange(128)
c.incr(2**power-1)
c.incr()
r = c.get()
```

what can go wrong?

- $\text{power} \leq 30$: $r = 2^{\text{power}}$
- $\text{power} = 31$: $r = 2^{-31}$: C overflow (silent wrap-around)
- $\text{power} \in [32, 62]$: Python OverflowError (overflow on int)
- $\text{power} \geq 63$: Python OverflowError (overflow on long)

Programs mixing C and Python

Python counter class in C

```
typedef struct {
    PyObject_HEAD;
    int counter;
} Counter;

static PyObject*
CounterIncr(Counter *self, PyObject *args) {
    int i = 1;
    if (!PyArg_ParseTuple(args, "|i", &i))
        return NULL;
    self->counter += i;
    Py_RETURN_NONE;
}

static PyObject* CounterGet(Counter *self) {
    return Py_BuildValue("i", self->counter);
}
```

Python client

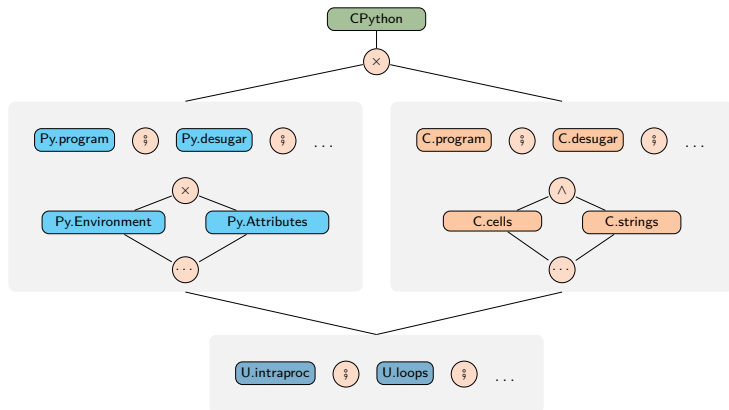
```
from counter import Counter
from random import randrange

c = Counter()
power = randrange(128)
c.incr(2**power-1)
c.incr()
r = c.get()
```

what can go wrong?

- $\text{power} \leq 30$: $r = 2^{\text{power}}$
- $\text{power} = 31$: $r = 2^{-31}$: C overflow (silent wrap-around)
- $\text{power} \in [32, 62]$: Python OverflowError (overflow on int)
- $\text{power} \geq 63$: Python OverflowError (overflow on long)

Domains for a C-Python value analysis



Analyze an AST containing **both** Python and C sources using

- Python domains: for Python code (some CPython API calls are translated to Python instructions)
- C domains: for C code (includes part of CPython implementation in C)
- **CPython domain**: translate objects between **shared C and Python heaps** with complementary **views** (boundary functions)

Application: small C libraries with Python bindings and unit tests (\simeq a few K lines)

[with R. Monat @ SAS'21]

Patch analysis for C

```
172 | 172 | /* Like fstatat, but cache the result. If ST->st_size is -1, the
173 | 173 | status has not been gotten yet. If less than -1, fstatat failed
174 | 174 | - with errno == -1 - ST->st_size. Otherwise, the status has already
175 | 175 | + with errno == ST->st_ino. Otherwise, the status has already
176 | 176 | been gotten, so return 0. */
177 | 177 | static int
178 | 178 | cache_fstatat (int fd, char const *file, struct stat *st, int flag)
179 | 179 | {
180 | 180 | - if (st->st_size == -1 && fstatat (fd, file, st, flag) != 0)
181 | 181 | + {
182 | 182 | + st->st_size = -2;
183 | 183 | + st->st_ino = errno;
184 | 184 | + }
185 | 185 | if (0 <= st->st_size)
186 | 186 | return 0;
187 | 187 | - errno = -1 - st->st_size;
188 | 188 | + errno = (int) st->st_ino;
189 | 189 | return -1;
190 | 190 | }
```

Diff from `remove.c` in Coreutils

- analyze a **pair** of programs
- whole-program **iteration** on both versions
- compute **semantic differences**: new program semantic and new domains
- goal: prove the **absence of regressions**

Application: 100-1500 line patches from Coreutils and Linux

[with D. Delmas @ SAS'19]

Exploitability analysis

Exploitable

```
void use(char * input) {
    char dest[10];
    strcpy(dest, input); // alarm!
}

void main() {
    char buf[100];
    fgets(buf, sizeof(buf), stdin);
    use(buf);
}
```

Non-exploitable

```
void use(char * input) {
    char dest[10];
    strcpy(dest, input); // alarm!
}

void main() {
    char buf[10]; // fixed!
    fgets(buf, sizeof(buf), stdin);
    use(rand() ? buf
           : "123456789012345");
}
```

- discover whether a runtime error **depends** on a **user-controllable value**
⇒ possibly **exploitable** bug
- **reduce alarm** number by reporting only exploitable ones
- combination of forward **value analysis** and **taint analysis**

[with F. Parolini @ VMCAI'24]

Finding witnesses for a bug

```
char buf[3];
int i = 0;
int j = input(0, 10); // input value?
int arr [3];
while (i < 5) {
    j += random(0, 1);
    i++;
}
buf[j] = '\0'; // bug
```

- given an **alarm**, find a **witness** proving the **presence** of the bug
- requires inferring **sufficient preconditions** for a given postcondition
 - infer user-**input** values
 - despite possible non-determinism (choices out of control of the user)
 - despite abstractions (approximations)
- the analysis proceeds **backwards** from the bug
- sufficient preconditions must be **under-approximated**

[with M. Milanese @ VMCAI'24]

Aspects of MOPSA to **improve**:

- **modular analyses** (beyond whole-program analyses)
 - analyze a function once, reuse the result many times
 - reuse **across different git projects** (libraries, files)
 - incremental analyses (add use context gradually)
 - develop more symbolic abstractions of the memory (separation properties)
- consider additional properties
 - further portability & impact analyses
 - **lightweight analyses**
- support **additional languages**, boost **multilanguage support**
 - work on OCaml support, OCaml-C bindings



<https://gitlab.com/mopsa/mopsa-analyzer>



The End
