

# Разработка алгоритмов статического поиска выходов за пределы динамического массива в C/C++ программах

И. Е. Громаковский

Научный руководитель: М. А. Лукин

Санкт-Петербургский национальный исследовательский университет  
информационных технологий, механики и оптики

- Статический поиск выходов за пределы динамического массива в C и C++.
- Работа с большими программами за разумное время.
- В общем случае задача поиска всех ошибок неразрешима.
- Находить как можно больше ошибок, минимизируя число ложных срабатываний.

- Программное обеспечение всегда было и остаётся подвержено ошибкам в коде.
- Выход за пределы массива — одна из главных уязвимостей с точки зрения безопасности.
- Особенно опасны ошибки в операционных системах и сетевых программах, для которых популярны языки C и C++.

- Меньше языковых конструкций, проще анализ.
- Анализируемый код ближе к реально выполняемому.
- Встроенные оптимизации и средства для анализа кода.
- Автоматически поддерживается любой язык, для которого есть компилятор в LLVM-IR.
- Static Single Assignment.

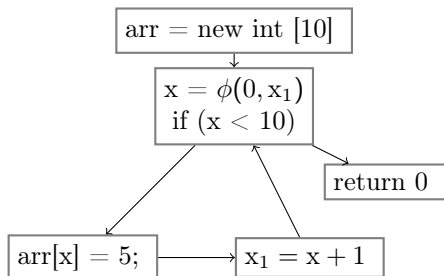
- Li L., Cifuentes C., Keynes N. Practical and Effective Symbolic Analysis for Buffer Overflow Detection.
- Symbolic ranges.
  - Symbolic expressions:  $15$ ,  $x + 1$ ,  $\perp$ ,  $\top$ .
  - Частичный порядок.
  - Symbolic range:  $[s1, s2]$ .
  - Операции:  $\cup, \cap, +, -, \times, \div$ .
- Define range,  $S_v$ :
  - в месте определения  $V$ .
- Use range,  $S_{v,P}$ :
  - в произвольном месте  $P$ .
- Выход за пределы массива размера  $n$  в инструкции  $P$ :
  - $S_{n,P}^{\max} \prec S_{\text{index},P}^{\max} \vee S_{\text{index},P}^{\min} \prec -1$ .

- Зависимости по данным.
  - Define range вычисляется через диапазоны аргументов.
  - $S_{(a+b)} = S_{a,P} + S_{b,P}, P = a + b$
  - $S_{(\phi(a,b))} = S_{a,P} \cup S_{b,P}, P = \phi(a, b)$
  - ...
- Зависимости потока управления.
  - Уточнение диапазона  $S_{V,P}$  на основании условных переходов на пути к  $P$ .
  - Условные переходы, связанные с  $V$ , такие что:
    - $P$  строго доминируется условным переходом;
    - $P$  достижима только из одного потомка условного перехода.

# Моноotonно изменяющиеся переменные

```
int main()
{
    int * arr = new int[10];
    for (int x = 0; x < 10; ++x)
        arr[x] = 5;

    return 0;
}
```

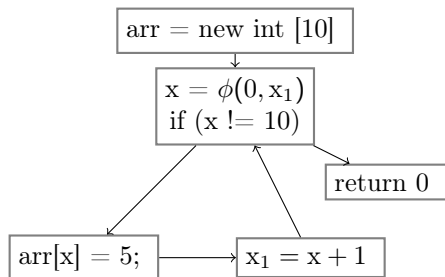


- Нет предикатов, ограничивающих значение снизу.
- Если  $x = \phi(a, b)$ ,  $b = f(x)$  и последовательность  $x, f(x), f(f(x)) \dots$  монотонна (например, возрастает), то применяется предикат  $x \geq a (x \leq a)$ .

# Обработка предиката «не равно»

```
int main()
{
    int * arr = new int[10];
    for (size_t x = 0; x != 10; ++x)
        arr[x] = 5;

    return 0;
}
```



- Алгоритм из статьи учитывает  $x \neq y$ , только если значение на границе отрезка.
- Если  $x = \phi(a, f(x))$ ,  $\exists c : y = f^c(a)$  и последовательность  $x, f(x), f(f(x)), \dots$  монотонна (например, возрастает), то применяется предикат  $x < y$ .



# Межпроцедурный анализ

```
#define ISDN_MAX_DRIVERS 32
#define ISDN_MAX_CHANNELS 64

struct isdn_driver * drivers[ISDN_MAX_DRIVERS];

struct isdn_driver * get_drv_by_nr(int di)
{
    struct isdn_driver * drv;
    drv = drivers[di];
    // ...
}

struct isdn_slot * get_slot_by_minor(int minor)
{
    int di, ch;
    struct isdn_driver * drv;
    for (di = 0; di < ISDN_MAX_CHANNELS; di++)
    {
        drv = get_drv_by_nr(di);
        // ..
    }
    // ..
}
```

- $32 \prec S_{di,P}^{\max}$ .
- $S_{di,P}^{\min} \prec -1$ .
- P — инструкция вызова.

- Триггер — условие ошибки:  $e_1 \prec e_2$ .
- Функции анализируются от вызываемой к вызывающей.
- Построение триггеров при анализе обращения к массиву.
- Проверка триггеров при анализе вызова функции.

# Сравнение: Multi Theft Auto

Multi Theft Auto 1.3.1: 760 тыс. строк кода, 182 мегабайта биткода.

—	TP	FP	Время работы (м:с)
Cppcheck	1	8	11:53
PVS-Studio	4	0	10:36
Splint	—	—	—
Li et al.	7	202	37:55
Улучшенная версия	7	4	42:49

TP (true positive) — число корректно найденных ошибок.

FP (false positive) — число ложных срабатываний.

## Сравнение: CMake

CMake 1.3.1: 350 тыс. строк кода, 31 мегабайт биткода.

—	TP	FP	Время работы (м:с)
Cppcheck	1	0	5:59
PVS-Studio	1	0	5:18
Splint	—	—	—
Li et al.	5	79	9:12
Улучшенная версия	5	3	11:03

TP (true positive) — число корректно найденных ошибок.

FP (false positive) — число ложных срабатываний.

# Сравнение: синтетические тесты

Синтетические тесты: различные использования массива, 12 ошибочных.

—	TP	FP	FN
Clang Analyzer	0	0	0
CppCheck	3	0	9
PVS-Studio	4	0	8
Splint	9	3	3
Li et al.	12	8	0
Улучшенная версия	12	0	0

TP (true positive) — число корректно найденных ошибок.

FP (false positive) — число ложных срабатываний.

FN (false negative) — число реальных ошибок, не найденных анализатором.

- За основу взят известный подход из статьи Li et al.
- Выявлены недостатки подхода.
- Сделаны улучшения, направленные на устранение выявленных недостатков.
- Получился анализатор, способный находить ошибки в больших программах, работающий лучше доступных анализаторов в некоторых случаях.

Спасибо за внимание!

Вопросы?