

WERYFIKACJA OPROGRAMOWANIA

SEMESTR ZIMOWY 2014/2015

Grzegorz Herman

Informatyka Analityczna
tcs@jagiellonian



WARUNKI ZALICZENIA

PUNKTY

- 5 zadań na Satori po 4 punkty każde
- egzamin: 3 z 4 pytań, po 6 punktów każde

OCENA

- z ćwiczeń: Satori
- końcowa: wszystko

PROGI

≤50%	2.0
50–60%	3.0
60–70%	3.5
70–80%	4.0
80–90%	4.5
>90%	5.0

BONUS

- implementacja metody z wykładu – podwyższenie oceny

PLAN ĆWICZEŃ

C/C++

- podstawowe narzędzia
- testy w modelu „black-box”
- analiza programów wielowątkowych

JAVA

- unit testing
- mock objects
- pokrycie kodu, mutation testing
- język specyfikacji JML
- instrumentacja

TESTOWANIE UI

- web-based UI
- desktop UI

PLAN WYKŁADÓW

WSTĘP

- grafowe reprezentacje programów

CZEŚĆ 1: ANALIZA DYNAMICZNA

- techniki instrumentacji
- wykrywanie data races
- analiza wpływu

CZEŚĆ 2: JAKOŚĆ I GENEROWANIE TESTÓW

- mutation testing
- testy pokrywające ścieżkę/punkt
- generowanie testów strukturalnych

PLAN WYKŁADÓW (CD.)

CZĘŚĆ 3: ANALIZA STATYCZNA

- wnioskowanie oparte o type inference
- analiza wskazywania
- przekroje

CZĘŚĆ 4: MODEL CHECKING

- logika Hoare'a
- logiki temporalne
- algorytmy model checking



- wskazówki
- optymalizacje

ANALIZA STATYCZNA

- nie uruchamia programu
- ogólne własności programu
- więcej informacji
- trudniejsza

ANALIZA DYNAMICZNA

- uruchamia program
- konkretny przebieg programu
- mniej informacji
- (względnie) prostsza



- feedback



- wskazówki
- optymalizacje

ANALIZA STATYCZNA

- nie uruchamia programu
- ogólne własności programu
- więcej informacji
- trudniejsza

ANALIZA DYNAMICZNA

- uruchamia program
- konkretny przebieg programu
- mniej informacji
- (względnie) prostsza



- feedback

WEJŚCIE: KOD ŹRÓDŁOWY

ZAŁOŻENIA

- język imperatywny
- pojedynczy wątek
- determinizm
- pojedyncza funkcja (analiza *intraproceduralna*)

UPROSZCZENIA WSTĘPNE

- dekonstrukcja struktur wysokiego poziomu
- przepływ sterowania zamieniony na skoki warunkowe

```
int ten() {  
    int i;  
    for (i=0; i<10; ++i);  
    return i;  
}
```

⇒

```
    i = 0;  
checkfor:  
    if (i>=10) goto endfor;  
    ++i;  
    goto checkfor;  
endfor:  
    return i;
```


WEJŚCIE: KOD ŹRÓDŁOWY

ZAŁOŻENIA

- język imperatywny
- pojedynczy wątek
- determinizm
- pojedyncza funkcja (analiza *intraproceduralna*)

UPROSZCZENIA WSTĘPNE

- dekonstrukcja struktur wysokiego poziomu
- przepływ sterowania zamieniony na skoki warunkowe

```
int ten() {  
    int i;  
    for (i=0; i<10; ++i);  
    return i;  
}
```

 \Rightarrow

```
    i = 0;  
checkfor:  
    if (i>=10) goto endfor;  
    ++i;  
    goto checkfor;  
endfor:  
    return i;
```

CONTROL FLOW GRAPH

BASIC BLOCK

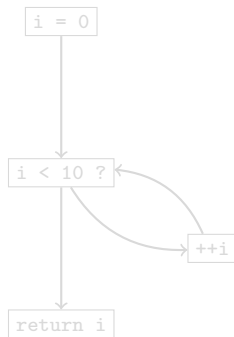
- liniowy ciąg instrukcji
- skoki „na zewnątrz” tylko z ostatniej
- skoki „do wewnątrz” tylko do pierwszej

CONTROL FLOW GRAPH $G = (V, E, s, t)$

- V – zbiór basic blocks
- $E \subseteq V \times V$ – możliwy przepływ sterowania
- $s \in V$ – instrukcja wejściowa
- $t \in V$ – instrukcja wyjściowa

REGULARYZACJA

- każdy $v \in V$ osiągalny z s
- t osiągalny z każdego $v \in V$



CONTROL FLOW GRAPH

BASIC BLOCK

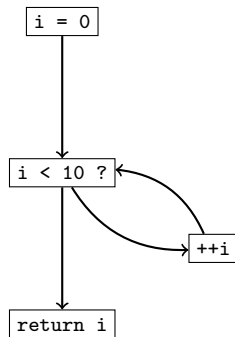
- liniowy ciąg instrukcji
- skoki „na zewnątrz” tylko z ostatniej
- skoki „do wewnątrz” tylko do pierwszej

CONTROL FLOW GRAPH $G = (V, E, s, t)$

- V – zbiór basic blocks
- $E \subseteq V \times V$ – możliwy przepływ sterowania
- $s \in V$ – instrukcja wejściowa
- $t \in V$ – instrukcja wyjściowa

REGULARYZACJA

- każdy $v \in V$ osiągalny z s
- t osiągalny z każdego $v \in V$



DOMINACJA

u DOMINUJE v

gdy każda ścieżka z s do v przechodzi przez u

u BEZPOŚREDNIO DOMINUJE v

gdy dodatkowo u nie dominuje żadnego innego dominatora v

BEZPOŚREDNIE DOMINATORY

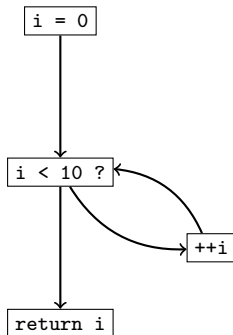
- tworzą drzewo o korzeniu w s
- można wyznaczyć w czasie prawie liniowym

$v \in \text{DOMINANCE FRONTIER}(u)$

- $u \rightsquigarrow w \rightarrow v$
- u dominuje w
- u nie dominuje v

POSTDOMINACJA

to dominacja od t po odwróconych krawędziach



DOMINACJA

u DOMINUJE v

gdy każda ścieżka z s do v przechodzi przez u

u BEZPOŚREDNIO DOMINUJE v

gdy dodatkowo u nie dominuje żadnego innego dominatora v

BEZPOŚREDNIE DOMINATORY

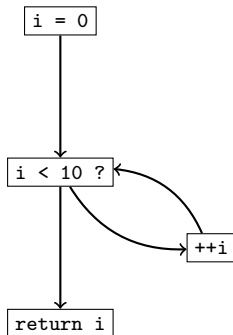
- tworzą drzewo o korzeniu w s
- można wyznaczyć w czasie prawie liniowym

$v \in \text{DOMINANCE FRONTIER}(u)$

- $u \rightsquigarrow w \rightarrow v$
- u dominuje w
- u nie dominuje v

POSTDOMINACJA

to dominacja od t po odwróconych krawędziach



DOMINACJA

 u DOMINUJE v

gdy każda ścieżka z s do v przechodzi przez u

 u BEZPOŚREDNIO DOMINUJE v

gdy dodatkowo u nie dominuje żadnego innego dominatora v

BEZPOŚREDNIE DOMINATORY

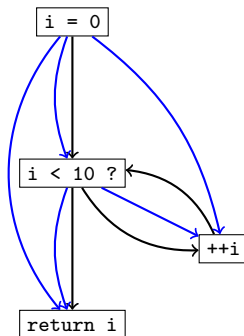
- tworzą drzewo o korzeniu w s
- można wyznaczyć w czasie prawie liniowym

 $v \in \text{DOMINANCE FRONTIER}(u)$

- $u \rightsquigarrow w \rightarrow v$
- u dominuje w
- u nie dominuje v

POSTDOMINACJA

to dominacja od t po odwróconych krawędziach



DOMINACJA

u DOMINUJE v

gdy każda ścieżka z s do v przechodzi przez u

u BEZPOŚREDNIO DOMINUJE v

gdy dodatkowo u nie dominuje żadnego innego dominatora v

BEZPOŚREDNIE DOMINATORY

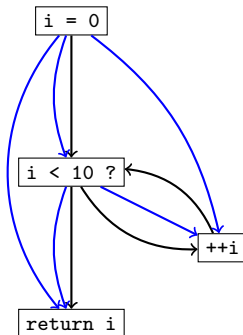
- tworzą drzewo o korzeniu w s
- można wyznaczyć w czasie prawie liniowym

$v \in \text{DOMINANCE FRONTIER}(u)$

- $u \rightsquigarrow w \rightarrow v$
- u dominuje w
- u nie dominuje v

POSTDOMINACJA

to dominacja od t po odwróconych krawędziach



DOMINACJA

u DOMINUJE v

gdy każda ścieżka z s do v przechodzi przez u

u BEZPOŚREDNIO DOMINUJE v

gdy dodatkowo u nie dominuje żadnego innego dominatora v

BEZPOŚREDNIE DOMINATORY

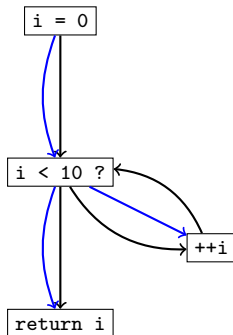
- tworzą drzewo o korzeniu w s
- można wyznaczyć w czasie prawie liniowym

$v \in \text{DOMINANCE FRONTIER}(u)$

- $u \rightsquigarrow w \rightarrow v$
- u dominuje w
- u nie dominuje v

POSTDOMINACJA

to dominacja od t po odwróconych krawędziach



DOMINACJA

u DOMINUJE v

gdy każda ścieżka z s do v przechodzi przez u

u BEZPOŚREDNIO DOMINUJE v

gdy dodatkowo u nie dominuje żadnego innego dominatora v

BEZPOŚREDNIE DOMINATORY

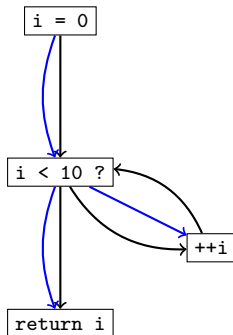
- tworzą drzewo o korzeniu w s
- można wyznaczyć w czasie prawie liniowym

$v \in \text{DOMINANCE FRONTIER}(u)$

- $u \rightsquigarrow w \rightarrow v$
- u dominuje w
- u nie dominuje v

POSTDOMINACJA

to dominacja od t po odwróconych krawędziach



DOMINACJA

u DOMINUJE v

gdy każda ścieżka z s do v przechodzi przez u

u BEZPOŚREDNIO DOMINUJE v

gdy dodatkowo u nie dominuje żadnego innego dominatora v

BEZPOŚREDNIE DOMINATORY

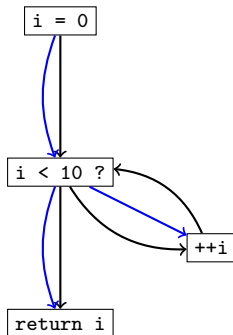
- tworzą drzewo o korzeniu w s
- można wyznaczyć w czasie prawie liniowym

$v \in \text{DOMINANCE FRONTIER}(u)$

- $u \rightsquigarrow w \rightarrow v$
- u dominuje w
- u nie dominuje v

POSTDOMINACJA

to dominacja od t po odwróconych krawędziach



DOMINACJA

u DOMINUJE v

gdy każda ścieżka z s do v przechodzi przez u

u BEZPOŚREDNIO DOMINUJE v

gdy dodatkowo u nie dominuje żadnego innego dominatora v

BEZPOŚREDNIE DOMINATORY

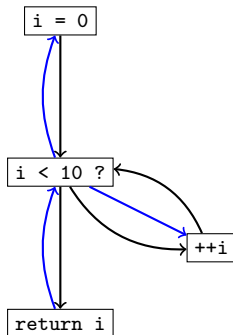
- tworzą drzewo o korzeniu w s
- można wyznaczyć w czasie prawie liniowym

$v \in \text{DOMINANCE FRONTIER}(u)$

- $u \rightsquigarrow w \rightarrow v$
- u dominuje w
- u nie dominuje v

POSTDOMINACJA

to dominacja od t po odwróconych krawędziach



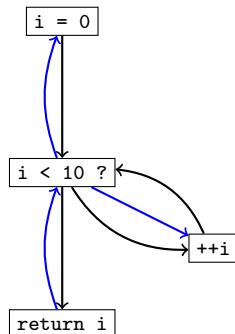
CONTROL DEPENDENCE GRAPH

 v CONTROL-DEPENDS ON u

- $\exists u \rightarrow w \rightsquigarrow v$
- v post-dominuje w (lub $v = w$)
- v nie post-dominuje u

INTUICJA

- u ma przynajmniej 2 wyjścia
- jedno z nich zawsze prowadzi do v
- drugie nie



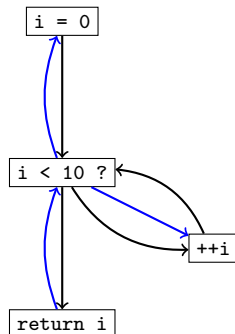
CONTROL DEPENDENCE GRAPH

 v CONTROL-DEPENDS ON u

- $\exists u \rightarrow w \rightsquigarrow v$
- v post-dominuje w (lub $v = w$)
- v nie post-dominuje u

INTUICJA

- u ma przynajmniej 2 wyjścia
- jedno z nich zawsze prowadzi do v
- drugie nie



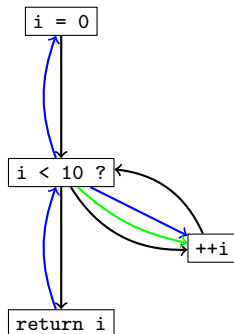
CONTROL DEPENDENCE GRAPH

 v CONTROL-DEPENDS ON u

- $\exists u \rightarrow w \rightsquigarrow v$
- v post-dominuje w (lub $v = w$)
- v nie post-dominuje u

INTUICJA

- u ma przynajmniej 2 wyjścia
- jedno z nich zawsze prowadzi do v
- drugie nie



DATA FLOW (DEPENDENCE) GRAPH

DEFINICJA ZMIENNEJ x

to instrukcja ustawiająca x

UŻYCIE ZMIENNEJ x

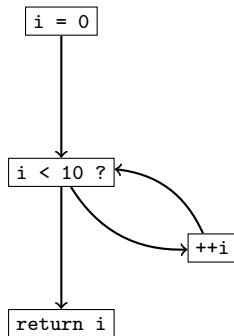
to instrukcja odczytująca x

ŚCIEŻKA WOLNA DLA x

to ścieżka w CFG omijająca definicje x

KRAWĘDŹ $u \rightarrow v$ w DFG

- u – definicja x
- v – użycie x
- $\exists u \rightsquigarrow v$ wolna dla x



DATA FLOW (DEPENDENCE) GRAPH

DEFINICJA ZMIENNEJ x

to instrukcja ustawiająca x

UŻYCIE ZMIENNEJ x

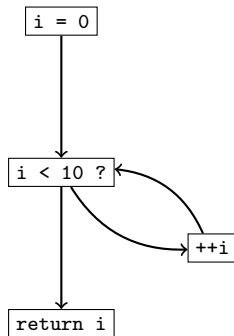
to instrukcja odczytująca x

ŚCIEŻKA WOLNA DLA x

to ścieżka w CFG omijająca definicje x

KRAWĘDŹ $u \rightarrow v$ w DFG

- u – definicja x
- v – użycie x
- $\exists u \rightsquigarrow v$ wolna dla x



DATA FLOW (DEPENDENCE) GRAPH

DEFINICJA ZMIENNEJ x

to instrukcja ustawiająca x

UŻYCIE ZMIENNEJ x

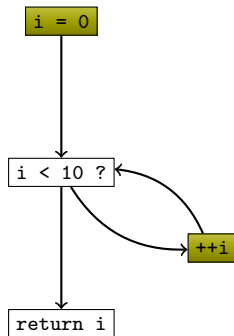
to instrukcja odczytująca x

ŚCIEŻKA WOLNA DLA x

to ścieżka w CFG omijająca definicje x

KRAWĘDŹ $u \rightarrow v$ w DFG

- u – definicja x
- v – użycie x
- $\exists u \rightsquigarrow v$ wolna dla x



DATA FLOW (DEPENDENCE) GRAPH

DEFINICJA ZMIENNEJ x

to instrukcja ustawiająca x

UŻYCIE ZMIENNEJ x

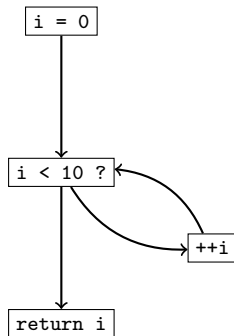
to instrukcja odczytująca x

ŚCIEŻKA WOLNA DLA x

to ścieżka w CFG omijająca definicje x

KRAWĘDŹ $u \rightarrow v$ w DFG

- u – definicja x
- v – użycie x
- $\exists u \rightsquigarrow v$ wolna dla x



DATA FLOW (DEPENDENCE) GRAPH

DEFINICJA ZMIENNEJ x

to instrukcja ustawiająca x

UŻYCIE ZMIENNEJ x

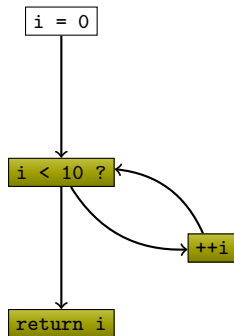
to instrukcja odczytująca x

ŚCIEŻKA WOLNA DLA x

to ścieżka w CFG omijająca definicje x

KRAWĘDŹ $u \rightarrow v$ w DFG

- u – definicja x
- v – użycie x
- $\exists u \rightsquigarrow v$ wolna dla x



DATA FLOW (DEPENDENCE) GRAPH

DEFINICJA ZMIENNEJ x

to instrukcja ustawiająca x

UŻYCIE ZMIENNEJ x

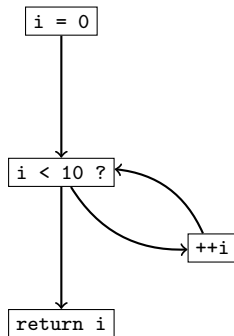
to instrukcja odczytująca x

ŚCIEŻKA WOLNA DLA x

to ścieżka w CFG omijająca definicje x

KRAWĘDŹ $u \rightarrow v$ w DFG

- u – definicja x
- v – użycie x
- $\exists u \rightsquigarrow v$ wolna dla x



DATA FLOW (DEPENDENCE) GRAPH

DEFINICJA ZMIENNEJ x

to instrukcja ustawiająca x

UŻYCIE ZMIENNEJ x

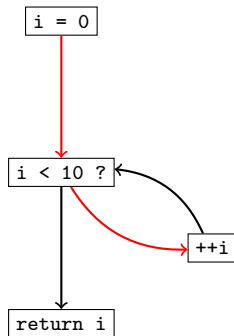
to instrukcja odczytująca x

ŚCIEŻKA WOLNA DLA x

to ścieżka w CFG omijająca definicje x

KRAWĘDŹ $u \rightarrow v$ w DFG

- u – definicja x
- v – użycie x
- $\exists u \rightsquigarrow v$ wolna dla x



DATA FLOW (DEPENDENCE) GRAPH

DEFINICJA ZMIENNEJ x

to instrukcja ustawiająca x

UŻYCIE ZMIENNEJ x

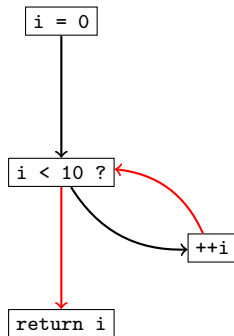
to instrukcja odczytująca x

ŚCIEŻKA WOLNA DLA x

to ścieżka w CFG omijająca definicje x

KRAWĘDŹ $u \rightarrow v$ w DFG

- u – definicja x
- v – użycie x
- $\exists u \rightsquigarrow v$ wolna dla x



DATA FLOW (DEPENDENCE) GRAPH

DEFINICJA ZMIENNEJ x

to instrukcja ustawiająca x

UŻYCIE ZMIENNEJ x

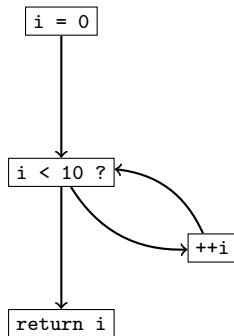
to instrukcja odczytująca x

ŚCIEŻKA WOLNA DLA x

to ścieżka w CFG omijająca definicje x

KRAWĘDŹ $u \rightarrow v$ w DFG

- u – definicja x
- v – użycie x
- $\exists u \rightsquigarrow v$ wolna dla x



DATA FLOW (DEPENDENCE) GRAPH

DEFINICJA ZMIENNEJ x

to instrukcja ustawiająca x

UŻYCIE ZMIENNEJ x

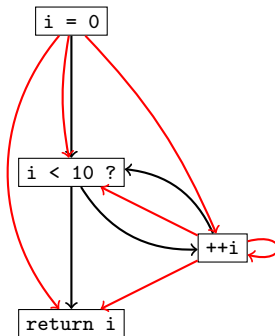
to instrukcja odczytująca x

ŚCIEŻKA WOLNA DLA x

to ścieżka w CFG omijająca definicje x

KRAWĘDŹ $u \rightarrow v$ w DFG

- u – definicja x
- v – użycie x
- $\exists u \rightsquigarrow v$ wolna dla x



DATA FLOW (DEPENDENCE) GRAPH

DEFINICJA ZMIENNEJ x

to instrukcja ustawiająca x

UŻYCIE ZMIENNEJ x

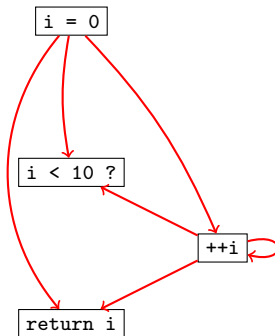
to instrukcja odczytująca x

ŚCIEŻKA WOLNA DLA x

to ścieżka w CFG omijająca definicje x

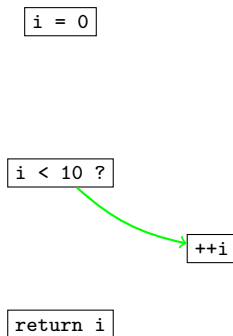
KRAWĘDŹ $u \rightarrow v$ w DFG

- u – definicja x
- v – użycie x
- $\exists u \rightsquigarrow v$ wolna dla x

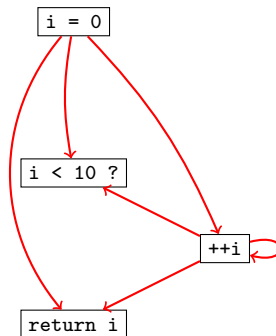


PROGRAM DEPENDENCE GRAPH

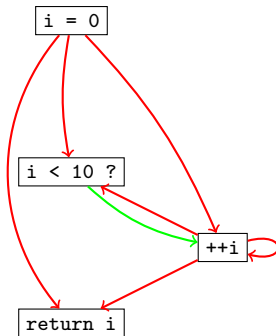
CONTROL DEPENDENCE



DATA DEPENDENCE



PROGRAM DEPENDENCE GRAPH



STATIC SINGLE ASSIGNMENT FORM

DEFINICJA ZMIENNEJ x

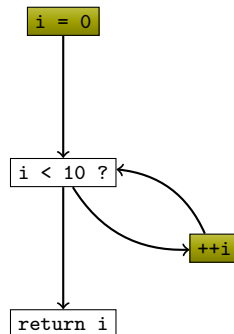
to instrukcja ustawiająca x

OGRANICZENIE SSA

każda zmienna ma dokładnie 1 definicję

KONSTRUKCJA

- nowa nazwa zmiennej w każdej definicji
- spotkanie 2+ definicji – sztuczna zmienna
- propagacja nowych nazw



STATIC SINGLE ASSIGNMENT FORM

DEFINICJA ZMIENNEJ x

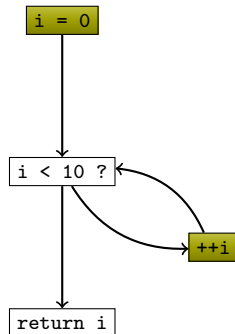
to instrukcja ustawiająca x

OGRANICZENIE SSA

każda zmienna ma dokładnie 1 definicję

KONSTRUKCJA

- nowa nazwa zmiennej w każdej definicji
- spotkanie 2+ definicji – sztuczna zmienna
- propagacja nowych nazw



STATIC SINGLE ASSIGNMENT FORM

DEFINICJA ZMIENNEJ x

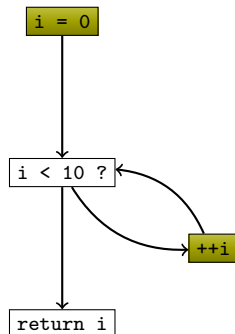
to instrukcja ustawiająca x

OGRANICZENIE SSA

każda zmienna ma dokładnie 1 definicję

KONSTRUKCJA

- nowa nazwa zmiennej w każdej definicji
- spotkanie 2+ definicji – sztuczna zmienna
- propagacja nowych nazw



STATIC SINGLE ASSIGNMENT FORM

DEFINICJA ZMIENNEJ x

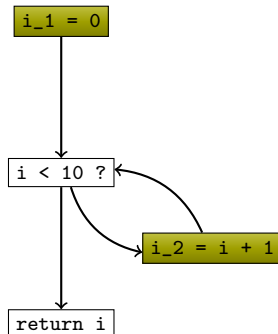
to instrukcja ustawiająca x

OGRANICZENIE SSA

każda zmienna ma dokładnie 1 definicję

KONSTRUKCJA

- nowa nazwa zmiennej w każdej definicji
- spotkanie 2+ definicji – sztuczna zmienna
- propagacja nowych nazw



STATIC SINGLE ASSIGNMENT FORM

DEFINICJA ZMIENNEJ x

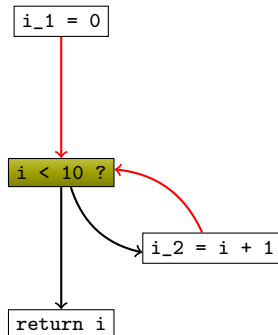
to instrukcja ustawiająca x

OGRANICZENIE SSA

każda zmienna ma dokładnie 1 definicję

KONSTRUKCJA

- nowa nazwa zmiennej w każdej definicji
- spotkanie 2+ definicji – sztuczna zmienna
- propagacja nowych nazw



STATIC SINGLE ASSIGNMENT FORM

DEFINICJA ZMIENNEJ x

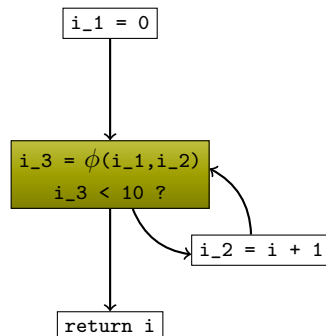
to instrukcja ustawiająca x

OGRANICZENIE SSA

każda zmienna ma dokładnie 1 definicję

KONSTRUKCJA

- nowa nazwa zmiennej w każdej definicji
- spotkanie 2+ definicji – sztuczna zmienna
- propagacja nowych nazw



STATIC SINGLE ASSIGNMENT FORM

DEFINICJA ZMIENNEJ x

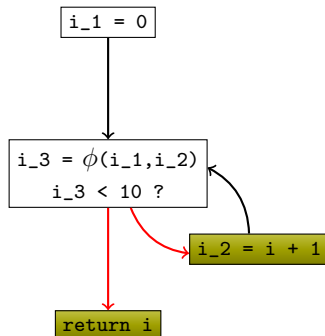
to instrukcja ustawiająca x

OGRANICZENIE SSA

każda zmienna ma dokładnie 1 definicję

KONSTRUKCJA

- nowa nazwa zmiennej w każdej definicji
- spotkanie 2+ definicji – sztuczna zmienna
- propagacja nowych nazw



STATIC SINGLE ASSIGNMENT FORM

DEFINICJA ZMIENNEJ x

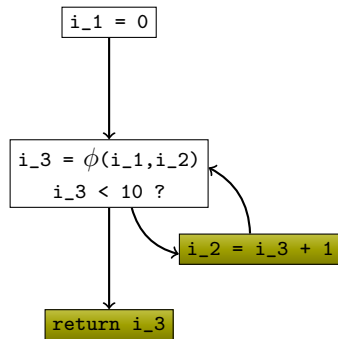
to instrukcja ustawiająca x

OGRANICZENIE SSA

każda zmienna ma dokładnie 1 definicję

KONSTRUKCJA

- nowa nazwa zmiennej w każdej definicji
- spotkanie 2+ definicji – sztuczna zmienna
- propagacja nowych nazw



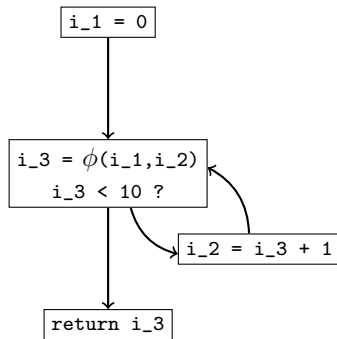
STATIC SINGLE INFORMATION FORM

OGRANICZENIA SSA

- każda zmienna ma dokładnie 1 definicję
- definicja zmiennej dominuje każde użycie
- każde użycie zmiennej post-dominuje definicję
- ϕ -użycia są w dominance frontier definicji
- σ -definicje są w postdominance frontier użyc

KONSTRUKCJA

- dla każdej niezależnej gałęzi – sztuczna zmienna
- propagacja nowych nazw



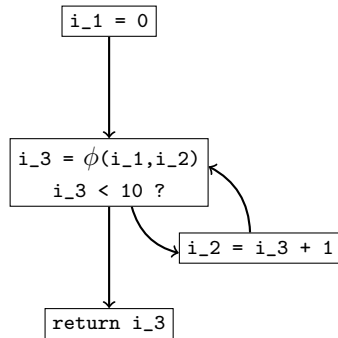
STATIC SINGLE INFORMATION FORM

OGRANICZENIA SSI

- każda zmienna ma dokładnie 1 definicję
- definicja zmiennej dominuje każde użycie
- każde użycie zmiennej post-dominuje definicję
- ϕ -użycia są w dominance frontier definicji
- σ -definicje są w postdominance frontier użyc

KONSTRUKCJA

- dla każdej niezależnej gałęzi – sztuczna zmienna
- propagacja nowych nazw



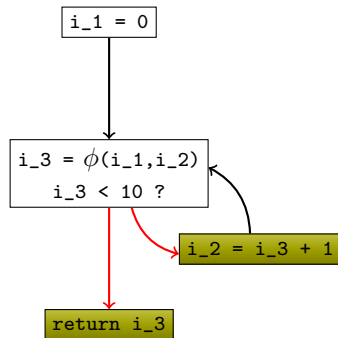
STATIC SINGLE INFORMATION FORM

OGRANICZENIA SSI

- każda zmienna ma dokładnie 1 definicję
- definicja zmiennej dominuje każde użycie
- każde użycie zmiennej post-dominuje definicję
- ϕ -użycia są w dominance frontier definicji
- σ -definicje są w postdominance frontier użyc

KONSTRUKCJA

- dla każdej niezależnej gałęzi – sztuczna zmienna
- propagacja nowych nazw



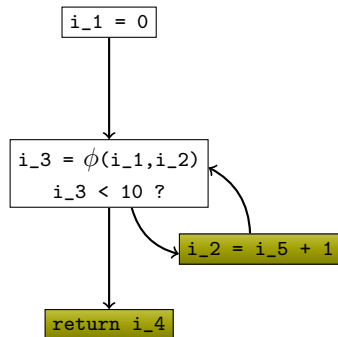
STATIC SINGLE INFORMATION FORM

OGRANICZENIA SSI

- każda zmienna ma dokładnie 1 definicję
- definicja zmiennej dominuje każde użycie
- każde użycie zmiennej post-dominuje definicję
- ϕ -użycia są w dominance frontier definicji
- σ -definicje są w postdominance frontier użyc

KONSTRUKCJA

- dla każdej niezależnej gałęzi – sztuczna zmienna
- propagacja nowych nazw



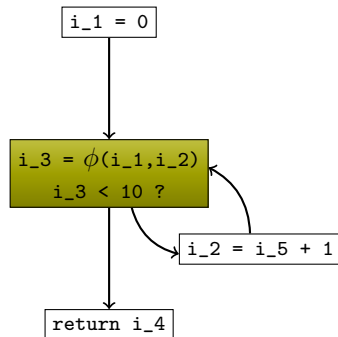
STATIC SINGLE INFORMATION FORM

OGRANICZENIA SSI

- każda zmienna ma dokładnie 1 definicję
- definicja zmiennej dominuje każde użycie
- każde użycie zmiennej post-dominuje definicję
- ϕ -użycia są w dominance frontier definicji
- σ -definicje są w postdominance frontier użyc

KONSTRUKCJA

- dla każdej niezależnej gałęzi – sztuczna zmienna
- propagacja nowych nazw



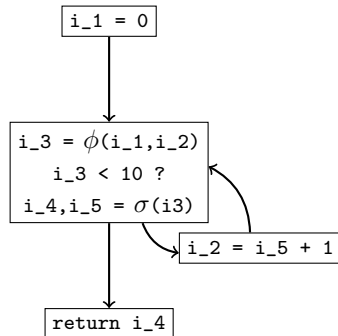
STATIC SINGLE INFORMATION FORM

OGRANICZENIA SSI

- każda zmienna ma dokładnie 1 definicję
- definicja zmiennej dominuje każde użycie
- każde użycie zmiennej post-dominuje definicję
- ϕ -użycia są w dominance frontier definicji
- σ -definicje są w postdominance frontier użyc

KONSTRUKCJA

- dla każdej niezależnej gałęzi – sztuczna zmienna
- propagacja nowych nazw



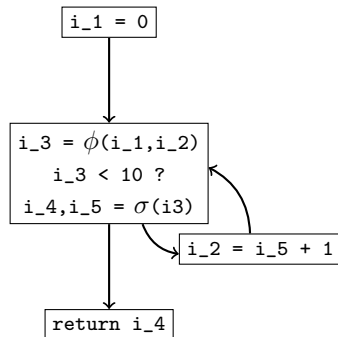
STATIC SINGLE INFORMATION FORM

OGRANICZENIA SSI

- każda zmienna ma dokładnie 1 definicję
- definicja zmiennej dominuje każde nie- ϕ -użycie
- każde użycie zmiennej post-dominuje nie- σ -definicję
- ϕ -użycia są w dominance frontier definicji
- σ -definicje są w postdominance frontier użyc

KONSTRUKCJA

- dla każdej niezależnej gałęzi – sztuczna zmienna
- propagacja nowych nazw



INSTRUMENTACJA – METODY

CZAS INSTRUMENTACJI

- statycznie w procesie kompilacji
- statycznie w momencie uruchamiania programu
- dynamicznie – w czasie działania programu

JĘZYK INSTRUMENTOWANY

- język źródłowy (wysokiego poziomu)
- język pośredni / bytecode
- kod binarny (zależny od architektury)

KOD INSTRUMENTACJI

- oddzielny proces, komunikacja przez IPC
- oddzielny wątek/wątki, komunikacja przez pamięć wspólną
- wbudowany, komunikacja przez wywołania funkcji
- wbudowany, inlining

INSTRUMENTACJA – METODY

INSTRUMENTACJA W CZASIE KOMPILACJI

- dostęp do informacji z analizy statycznej
- pełny inlining oraz inne optymalizacje
- wymagane wsparcia kompilatora/języka
- brak obsługi generowanego/modyfikowanego kodu

INSTRUMENTACJA – METODY

INSTRUMENTACJA STATYCZNA PRZY URUCHAMIANIU PROGRAMU

- niezależność od języka źródłowego
- możliwa analiza statyczna
- możliwe optymalizacje
- brak obsługi generowanego/modyfikowanego kodu
- trudności z odróżnieniem kodu od danych

INSTRUMENTACJA – METODY

INSTRUMENTACJA DYNAMICZNA

- niezależność od języka źródłowego
- łatwa identyfikacja kodu
- obsługa kodu generowanego/modyfikowanego
- dostępny jedynie lokalny widok programu
- lokalne optymalizacje możliwe, lecz trudne

INSTRUMENTACJA – METODY

INSTRUMENTACJA Z ODDZIELNEGO PROCESU

- pełna izolacja
- przezroczystość niemal darmowa
- duży koszt komunikacji
- brak dostępu do stanu programu

INSTRUMENTACJA – METODY

INSTRUMENTACJA Z DEDYKOWANEGO WĄTKU

- izolacja kodu
- przezroczystość względnie łatwa do osiągnięcia
- bezpośredni dostęp do stanu programu
- instrumentacja wymaga context-switch
- duży koszt synchronizacji

INSTRUMENTACJA – METODY

INSTRUMENTACJA Z „SHADOW THREADS”

- izolacja kodu
- przezroczystość względnie łatwa do osiągnięcia
- bezpośredni dostęp do stanu programu
- instrumentacja wymaga context-switch
- podwojone zużycie zasobów dla wątków

INSTRUMENTACJA – METODY

INSTRUMENTACJA WBUDOWANA W KOD

- minimalny koszt wydajnościowy
- bezpośredni dostęp do stanu programu
- brak izolacji
- przezroczystość trudna lub niemożliwa do osiągnięcia

METODY MODYFIKACJI KODU

- fault insertion
- jump insertion
- rekompilacja JIT

REKOMPILACJA DYNAMICZNA

SCHEMAT METODY

- basic block – liniowy fragment kodu (oryginalnego)
- oryginalny kod nie jest wykonywany bezpośrednio
- przy pierwszym wejściu do bloku – rekompilacja z dodaniem instrumentacji
- cache skompilowanych fragmentów

REKOMPILACJA DYNAMICZNA: CONTROL FLOW

RODZAJE SKOKÓW (CONTROL FLOW)

- bezpośrednie: adres docelowy stały dla instrukcji
- pośrednie: cel zależny od danych (w tym `return`)
- warunkowość skoku nie wpływa na jego bezpośredniość!

SKOKI BEZPOŚREDNIE

- zastąpione przez skok do odpowiedniego skompilowanego bloku

SKOKI POŚREDNIE

- wyszukiwanie adresu w skompilowanych blokach

REKOMPILACJA DYNAMICZNA: CONTROL FLOW

WYWOŁANIA PROCEDUR

- rzeczywiste adresy zmienione
- bezpośrednie wykorzystanie `call` i `ret` – stos zawiera inne wartości, ważne szczególnie przy PIC
- kodowanie jako `push/pop + jmp`
- bardziej złożona obsługa błędów

REKOMPILACJA DYNAMICZNA: CONTROL FLOW

OPTYMALIZACJA CZĘSTYCH ŚCIEŻEK WYKONANIA

- trace – ciąg bloków wykonywanych kolejno
- budowane na „ciepłych” blokach
- tworzone zachłannie aż do istniejącego trace’u, bezpośredniego skoku wstecznego lub limitu długości
- skoki pośrednie zastąpione przez porównanie i skok bezpośredni

REKOMPILACJA DYNAMICZNA: REJESTRY

PODEJŚCIE PODSTAWOWE

- wykorzystujemy oryginalne rejestry
- konieczność zapisania/odzyskania stanu przed/po instrumentacji

OPTYMALIZACJA

- algorytm realokacji rejestrów
- register renaming, zapisany dla każdego bloku
- kod wyrównujący przy każdym skoku
- analiza żywotności (szczególnie ważna przy `eflags`)

REKOMPILACJA DYNAMICZNA: STOS

WSPÓLNY STOS

- program może sięgać poza szczyt
- ryzyko przepełnienia przy instrumentacji

ODDZIELNY STOS

- konieczność przełączania
- utrata sprzętowej predykcji adresu powrotu

REKOMPILACJA DYNAMICZNA: BIBLIOTEKI

OBSŁUGA BIBLIOTEK

- mapa pamięci dla rozróżnienia kodu programu i bibliotek
- możliwe różne polityki instrumentacji

PROBLEMY/NIEBEZPIECZEŃSTWA

- callbacks – adresy kodu przekazywane do biblioteki
- współdzielenie biblioteki przez instrumentację i program

REKOMPILACJA DYNAMICZNA: ZASOBY SYSTEMOWE

OGRANICZENIA ZASOBÓW

- wspólne
- możliwa konieczność ukrywania zużycia

IZOLACJA ZASOBÓW

- ukrywanie
- możliwe niepożądane interakcje (np. `sync`)
- file descriptors – numer to też zasób!

REKOMPILACJA DYNAMICZNA: WIELOWĄTKOWOŚĆ

DOSTĘP DO DANYCH INSTRUMENTACJI

- adres w rejestrze
- stały adres w pamięci thread-local
- stały adres globalny – skompilowane bloki muszą być thread-local

SYNCHRONIZACJA

- wspólna synchronizacja programu i instrumentacji – niebezpieczeństwo deadlock'u
- gruboziarnista synchronizacja instrumentacji – ograniczenie wydajności
- w przeciwnym wypadku instrumentacja musi być re-entrant

REKOMPILACJA DYNAMICZNA: ZMIANY KODU

KOD GENEROWANY

- obsługiwany identycznie jak „zwykły”

MODYFIKACJE KODU

- odpowiedni blok wyrzucany z cache'u
- konieczność usunięcia wchodzących skoków bezpośrednich
- blok modyfikujący sam siebie (x86) – sprawdzanie spójności po każdej instrukcji
- modyfikacja pomiędzy wątkami (x86) – przy 3+ wątkach wymaga instrukcji synchronizującej, na której można podzielić blok; dla 2 wątków nierozwiązane (?)

REKOMPILACJA DYNAMICZNA: OBSŁUGA BŁĘDÓW

BŁĘDY PODCZAS INSTRUMENTACJI

- izolowane od aplikacji
- co robić, gdy nie można kontynuować instrumentacji?

BŁĘDY W REKOMPILOWANYM KODZIE

- wykryte przy kompilacji (np. invalid opcode) – podział bloku przed błędem
- wykryte w czasie działania – konieczność przetłumaczenia kontekstu

REKOMPILACJA DYNAMICZNA: DEBUGGER

„OBSŁUGA” DEBUGGERA

- pełna izolacja nieosiągalna
- komplikacje przy tłumaczeniu kontekstu
- wątek debuggera – instrumentowany czy nie?

ANALIZA WPŁYWU (TAINT ANALYSIS)

ZARYS TECHNIKI

- podział źródeł danych na „bezpieczne” i „niebezpieczne” (znaczone)
- śledzenie przepływu znaczone (tainted) danych
- analiza interesujących sytuacji

ZASTOSOWANIA

- wykrywanie ataków typu code-injection, SQL injection, cross-site-scripting
- generowanie filtrów bezpieczeństwa
- analiza przepływu zastrzeżonych informacji
- analiza oprogramowania typu malware
- wykrywanie naruszeń protokołów
- ...

TAINT ANALYSIS

MIARY WPŁYWU

- jakościowa (true/false)
- z uwzględnieniem źródła
- ilościowe: liczba bitów, entropia/channel capacity
- przybliżone: przedziały, zbiory wartości

TYPOWO WYKRYWANE SYTUACJE

- wykonanie znaczonego kodu
- skok pod znaczony adres
- odczyt/zapis pod znaczony adres

TAINT ANALYSIS – ŹRÓDŁA BŁĘDÓW

OVERTAINTING (FALSE POSITIVES)

- zbyt mała szczegółowość
- propagacja przez obecne w programie warunki bezpieczeństwa

UNDERTAINTING (FALSE NEGATIVES)

- implicit flow

TAINT ANALYSIS – WYZWANIA

ZNAKOWANIE WSKAŹNIKÓW

- uwzględnione: overtainting, zwłaszcza przy strukturach danych o bogatej strukturze pamięciowej
- pominięte: undertainting (ilościowo niewielki, ale istotny dla wykrywania buffer overruns itp.)

CONTROL FLOW TAINT

- jedynie analiza dynamiczna: niemożliwy do wykrycia
- analiza statyczna (preprocessing): overtainting

CONTROL FLOW TAINT

OBRANIE KONKRETNEJ ŚCIEŻKI WYKONANIA ZAWIERA INFORMACJE

```
if (x==true) y = true; else y = false;  
// x==y
```

NIEOBRANIE KONKRETNEJ ŚCIEŻKI WYKONANIA RÓWNIEŻ

```
y = z = false;  
if (x==false) z = true;  
if (z==false) y = true;  
// x==y
```

ROZWIĄZANIE KONSERWATYWNE

wszystkie zmienne przypisywane we **wszystkich** instrukcjach zależnych (control-flow dependent) od znakowanej zmiennej muszą być znakowane.

TAINT ANALYSIS – WYZWANIA

USUWANIE OZNACZENIA (SANITIZATION)

- funkcje stałe (np. `xor eax, eax`) i jednokierunkowe
- konstrukcje kompilacji (np. `switch`)
- rozwiązanie: adnotacje (nie zawsze możliwe)

OPÓŹNIONE WYKRYCIE

- przykład: `return address overwrite` – wykryty dopiero w momencie skoku
- rozwiązanie: postprocessing logów wykonania (nie zawsze praktyczne)