

# Abstract Semantic Differencing for Numerical Programs

**Nimrod Partush**

Eran Yahav

Technion, Israel

# Semantic differencing



Characterize **semantic difference** between **similar programs**

# Motivating example

```
1. if (input % 2 == 0) goto 2 else goto 4
2. s := input+2
3. goto 5
4. s := input+3
5.
6. ptr := realloc(ptr,s)
7. // use ptr[0], ptr[1], ... ptr[input-1]
```

$$2^{32} - 3 \leq \text{input} \leq 2^{32} - 1$$

```
1. if (input % 2 == 0) goto 2
2. s := input+2
3. goto 5
4. s := input+3
5. + if (s>input) goto 6 else goto ERROR
6. ptr := realloc(ptr,s)
7. // use ptr[0], ptr[1], ... ptr[input-1]
```

# Abstract semantic differencing

- Use **abstract interpretation** to **prove equivalence** between two program versions
- Or **characterize their difference**
  - find (an abstraction of) **all inputs** that lead to **different output**
- **Sound**
  - never miss a difference
- **Precise**
  - report few false differences

# Equivalence under abstraction

```
int sign(int x) {
  int sgn;
  if (x < 0)
    sgn = -1
  else
    sgn = 1

  return sgn
}
```

$$\text{sign}(x) = \begin{cases} -1 & , x < 0 \\ 1 & , x \geq 0 \end{cases}$$

$\text{sgn} \mapsto [-1, -1]$

$\text{sgn} \mapsto [1, 1]$

$\text{sgn} \mapsto [-1, -1] \sqcup [1, 1]$

$\text{sgn} \mapsto [-1, 1]$

```
int sign(int x) {
  int sgn;
  if (x < 0)
    sgn = -1
  else
    sgn = 1
  + if (x == 0)
  +   sgn = 0
  return sgn
}
```

$$\text{sign}(x) = \begin{cases} 1 & , x > 0 \\ 0 & , x = 0 \\ -1 & , x < 0 \end{cases}$$

$\text{sgn} \mapsto [-1, -1]$

$\text{sgn} \mapsto [1, 1]$

$\text{sgn} \mapsto [-1, -1] \sqcup [1, 1]$

$\text{sgn} \mapsto [0, 0]$

$\text{sgn} \mapsto [0, 0] \sqcup [-1, 1]$

$\text{sgn} \mapsto [-1, 1]$

Equivalence under abstraction does not entail  
equivalence between the concrete values it represents

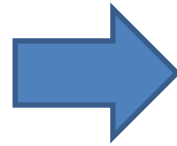
# Our approach

- Create a **correlating program**  $P \bowtie P'$  which captures behaviors of  $P$  and  $P'$
- Analyze  $P \bowtie P'$  using a partially disjunctive **correlating abstract domain**
  - Track equivalence between variables of  $P$  and  $P'$
  - Join states with the same equivalence relation (**partitioning**)

# Correlating Program $P \bowtie P'$

- We create a new syntactic object that combines  $P$  and  $P'$

```
int sign(int x) {  
  int sgn;  
  if (x < 0)  
    sgn = -1  
  else  
    sgn = 1  
  
  return sgn  
}  
  
int sign'(int x') {  
  int sgn';  
  if (x' < 0)  
    sgn' = -1  
  else  
    sgn' = 1  
  + if (x' == 0)  
  +   sgn' = 0  
  return sgn'  
}
```

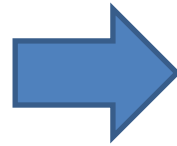


```
int sign  $\bowtie$  sign'(int x) {  
  int x' = x;  
  int sgn, sgn' = sgn;  
  (x < 0)  $\rightarrow$  sgn = -1  
  (x' < 0)  $\rightarrow$  sgn' = -1  
  (x  $\geq$  0)  $\rightarrow$  sgn = 1  
  (x'  $\geq$  0)  $\rightarrow$  sgn' = 1  
  (x' == 0)  $\rightarrow$  sgn' = 0  
}
```

# Correlating Program $P \bowtie P'$

- We create a new syntactic object that combines  $P$  and  $P'$

```
int sign(int x) {  
  int sgn;  
  if (x < 0)  
    sgn = -1  
  else  
    sgn = 1  
  
  return sgn  
}  
  
int sign'(int x') {  
  int sgn';  
  if (x' < 0)  
    sgn' = -1  
  else  
    sgn' = 1  
  + if (x' == 0)  
  +   sgn' = 0  
  return sgn'  
}
```



```
int sign $\bowtie$ sign'(int x) {  
  int x' = x;  
  int sgn, sgn' = sgn;  
  guard g1 = (x < 0);  
  guard g1' = (x' < 0);  
  if (g1) sgn = -1;  
  if (g1') sgn' = -1;  
  if (!g1) sgn = 1;  
  if (!g1') sgn' = 1;  
  guard g2' = (x' == 0);  
  if (g2') sgn' = 0;  
  retval = sgn;  
  retval' = sgn';  
}
```



# Correlating abstract domain

- Maintain **direct correlation** between values in the programs P and P'
  - Use an **relational** abstraction that **captures equivalences**

$\{ x < 0, \text{sgn} \mapsto -1 \}$

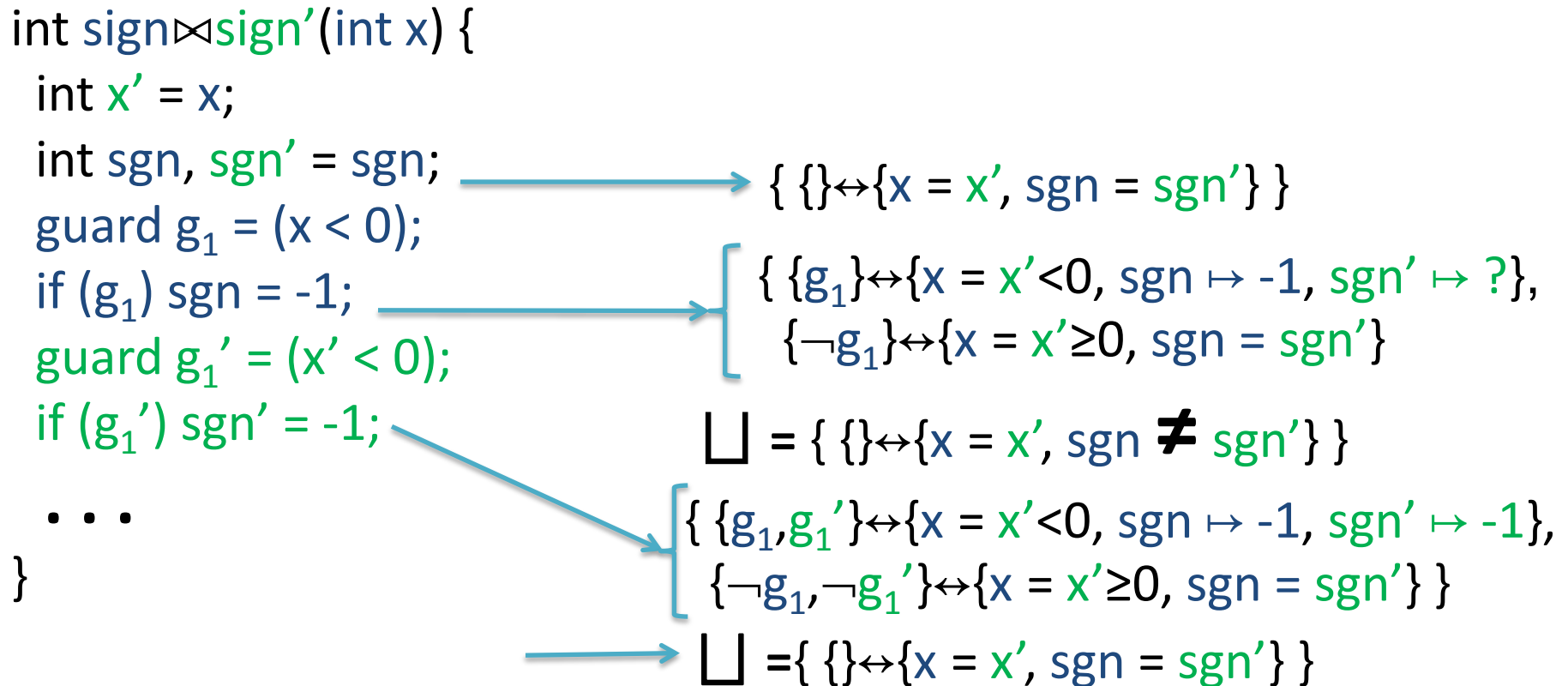
$\{ x' < 0, \text{sgn}' \mapsto -1 \}$



$\{ g_1, g_1' \} \leftrightarrow \{ x = x' < 0, \text{sgn} = \text{sgn}' \mapsto -1 \}$

- we use a **partially disjunctive domain** since we need to **delay joining**

# Delay $\sqcup$ to preserve equivalence



# Delay $\sqcup$ to preserve equivalence

- We want to join at locations in  $P \bowtie P'$  where equivalence is more likely to hold
  - After “matched” instructions both ran
- We call these **correlation points**
  - Part of  $P \bowtie P'$  creation process

# Picking correlation points in $P \bowtie P'$

```
int retval;  
int sign(int x) {  
  int sgn;  
  guard g1 = (x < 0);  
  if (g1) sgn = -1;  
  if (!g1) sgn = 1;  
  
  retval = sgn;  
}
```

```
int retval';  
int sign(int x') {  
  int sgn';  
  guard g1' = (x' < 0);  
  if (g1') sgn' = -1;  
  if (!g1') sgn' = 1;  
  
  guard g2' = (x' == 0);  
  if (g2') sgn' = 0;  
  retval' = sgn';  
}
```

```
int retval, retval';  
int sign $\bowtie$ sign'(int x) {  
  int x' = x;  
  int sgn, sgn' = sgn;  
  guard g1 = (x < 0);  
  if (g1) sgn = -1;  
  guard g1' = (x' < 0);  
  if (g1') sgn' = -1;  
  if (!g1) sgn = 1;  
  if (!g1') sgn' = 1;  
  guard g2' = (x' == 0);  
  if (g2') sgn' = 0;  
  retval = sgn;  
  retval' = sgn';  
}
```

# Correlating Program $P \bowtie P'$

- $P \bowtie P'$  is a **reduction over  $P \times P'$**  that's better for tracking equivalences and finding differences
  - **Construct the program in a way that matches the abstraction**
  - **brings matched instructions closer together**
- In general, can search the space of potential correlating programs
  - As well as correlation points

# Delay $\sqcup$ to preserve equivalence

```
int sign  $\bowtie$  sign'(int x) {
```

```
  int x' = x;
```

```
  int sgn, sgn' = sgn;
```

```
  guard g1 = (x < 0);
```

```
  if (g1) sgn = -1;
```

```
  guard g1' = (x' < 0);
```

```
  if (g1') sgn' = -1;
```

```
  if (!g1) sgn = 1;
```

```
  if (!g1') sgn' = 1;
```

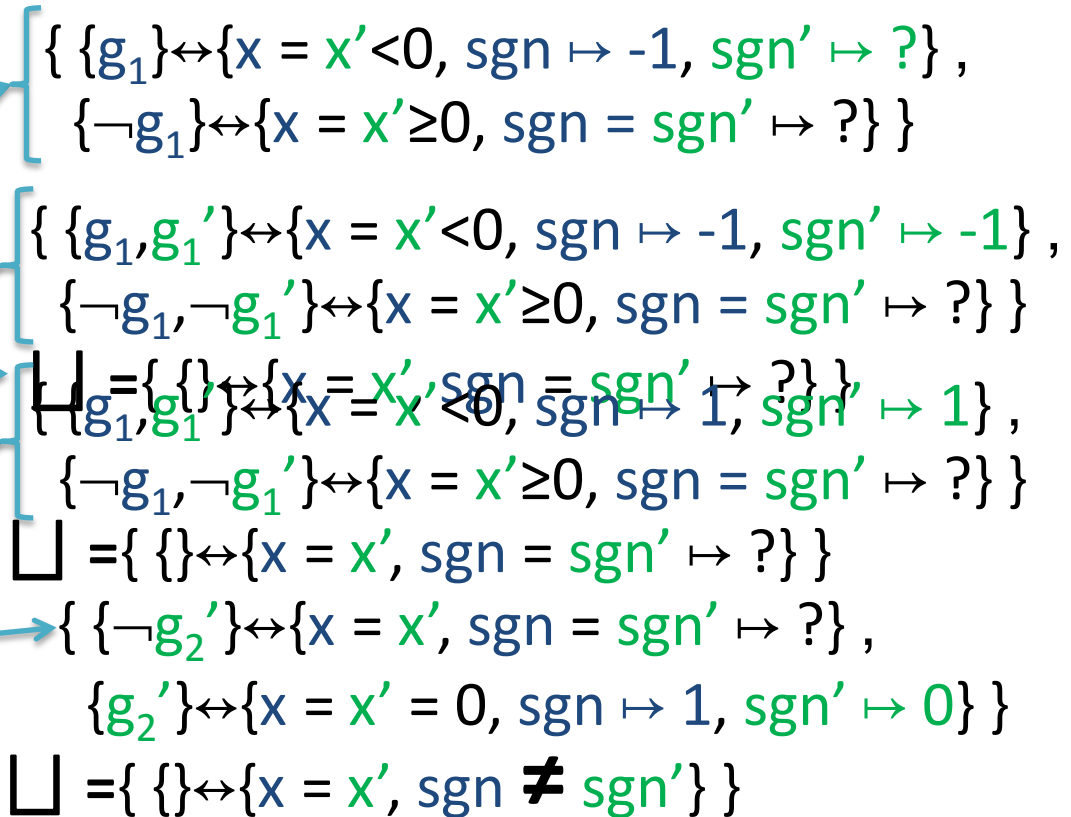
```
  guard g2' = (x' == 0);
```

```
  if (g2') sgn' = 0;
```

```
  retval = sgn;
```

```
  retval' = sgn';
```

```
}
```



# Partitioning based on equivalence

- Join abstract states based on the **equivalences they preserve**
  - the set of variables that hold equivalence
  - disjunction size bound at  $2^{|\text{VAR}|}$
  - lose some information, but maintain what's important (equivalence)

# Correlating Analysis for $P \bowtie P'$

```
int sign  $\bowtie$  sign'(int x) {  
  int x' = x;  
  int sgn, sgn' = sgn;  
  guard g1 = (x < 0);  
  if (g1) sgn = -1;  
  guard g1' = (x' < 0);  
  if (g1') sgn' = -1;  
  if (!g1) sgn = 1;  
  if (!g1') sgn' = 1;  
  guard g2' = (x' == 0);  
  if (g2') sgn' = 0;  
  retval = sgn;  
  retval' = sgn';  
}
```

$\{ \neg g_2' \} \leftrightarrow \{ x = x', \text{sgn} = \text{sgn}' \mapsto ? \},$

$\{ g_2' \} \leftrightarrow \{ x = x' = 0, \text{sgn} \mapsto 1, \text{sgn}' \mapsto 0 \}$

eqv

$\sqcup_{\text{eqv}} = \{ \neg g_2' \} \leftrightarrow \{ x = x', \text{sgn} = \text{sgn}' \mapsto ? \},$   
 $\{ g_2' \} \leftrightarrow \{ x = x' = 0, \text{sgn} \mapsto 1, \text{sgn}' \mapsto 0 \}$

diff



# Semantic differencing for **loops**

**P**

```
int foo(int x, int y, int z) {  
  while (x > 0) {  
    if (z > 0)  
      - y++;  
    x--;  
  }  
  return y;  
}
```

**P'**

```
int foo'(int x', int y', int z') {  
  while (x' > 0) {  
    if (z' > 0)  
      + y'--;  
    x'--;  
  }  
  return y';  
}
```

- These programs differ for cases where **z>0** and are otherwise equivalent

# Semantic differencing for loops

```
int foo  $\bowtie$  foo'(int x, int y, int z) {
```

```
  int x' = x, y' = y, z' = z;
```

```
  loop:
```

```
    guard g1 = (x > 0);
```

```
  loop':
```

```
    guard g'1 = (x' > 0);
```

```
    guard g2 = (z > 0);
```

```
    guard g'2 = (z' > 0);
```

```
    if (g1 && g2) y++;
```

```
    if (g'1 && g'2) y'--;
```

```
    if (g1) x--;
```

```
    if (g'1) x'--;
```

```
    if (g1) goto loop;
```

```
    if (g'1) goto loop';
```

```
}
```

$$\begin{aligned} & \{ \{g_1, g_2\} \leftrightarrow \{z = z' > 0, x = x' > 0, y = y' + 1\}, \\ & \{g_1, \neg g_2\} \leftrightarrow \{z = z' \leq 0, x = x' > 0, y = y'\}, \\ & \{\neg g_1\} \leftrightarrow \{z = z', x = x' \leq 0, y = y'\} \} \end{aligned}$$

$$\begin{aligned} & \{ \{g_1, g_1, g_2, g_2\} \leftrightarrow \{z = z' > 0, x = x' > 0, y = y' + 2\} \\ & \{g_1, g_1, \neg g_2, \neg g_2\} \leftrightarrow \{z = z' \leq 0, x = x' > 0, y = y'\} \\ & \{g_1, \neg g_1\} \leftrightarrow \{z = z', x = x', y = y'\} \\ & \{g_1, g_1\} \leftrightarrow \{z = z', x = x', y = y'\} \\ & \{g_1, g_1, g_2, g_2\} \leftrightarrow \{z = z' > 0, x = x' > 0, y = y' + 2\} \} \end{aligned}$$

# Semantic differencing for loops

```
int foo  $\bowtie$  foo'(int x, int y, int z) {
```

```
  int x' = x, y' = y, z' = z;
```

```
  loop:
```

```
    guard g1 = (x > 0);
```

```
  loop':
```

```
    guard g'1 = (x' > 0);
```

```
    guard g2 = (z > 0);
```

```
    guard g'2 = (z' > 0);
```

```
    if (g1 && g2) y++;
```

```
    if (g'1 && g'2) y'--;
```

```
    if (g1) x--;
```

```
    if (g'1) x'--;
```

```
    if (g1) goto loop;
```

```
    if (g'1) goto loop';
```

```
}
```

$\{ \{ \} \leftrightarrow \{ z = z', x = x', y = y' \} ,$   
 $\{ g_1, g_1, g_2, g_2 \} \leftrightarrow \{ z = z' > 0, x = x' > -1, y = y' + 2 \} \}$   
 $\{ \{ \} \leftrightarrow \{ z = z', x = x', y = y' \} ,$   
 $\{ g_1, g_1, g_2, g_2 \} \leftrightarrow \{ z = z' > 0, x = x' > -1, y = y' + 4 \} \}$

...

?

# Semantic differencing for **loops**

- We need to **widen in the powerset domain**, but **which sub-states** should be matched?
  - Strategy 1: Widen-by-equivalence
  - Strategy 2: Widen-by-guards

# Widen-by-equivalence

```
int foo  $\bowtie$  foo'(int x, int y, int z) {
```

```
  int x' = x, y' = y, z' = z;
```

```
  loop:
```

```
    guard g1 = (x > 0);
```

```
  loop':
```

```
    guard g'1 = (x' > 0);
```

```
    guard g2 = (z > 0);
```

```
    guard g'2 = (z' > 0);
```

```
    if (g1 && g2) y++;
```

```
    if (g'1 && g'2) y'--;
```

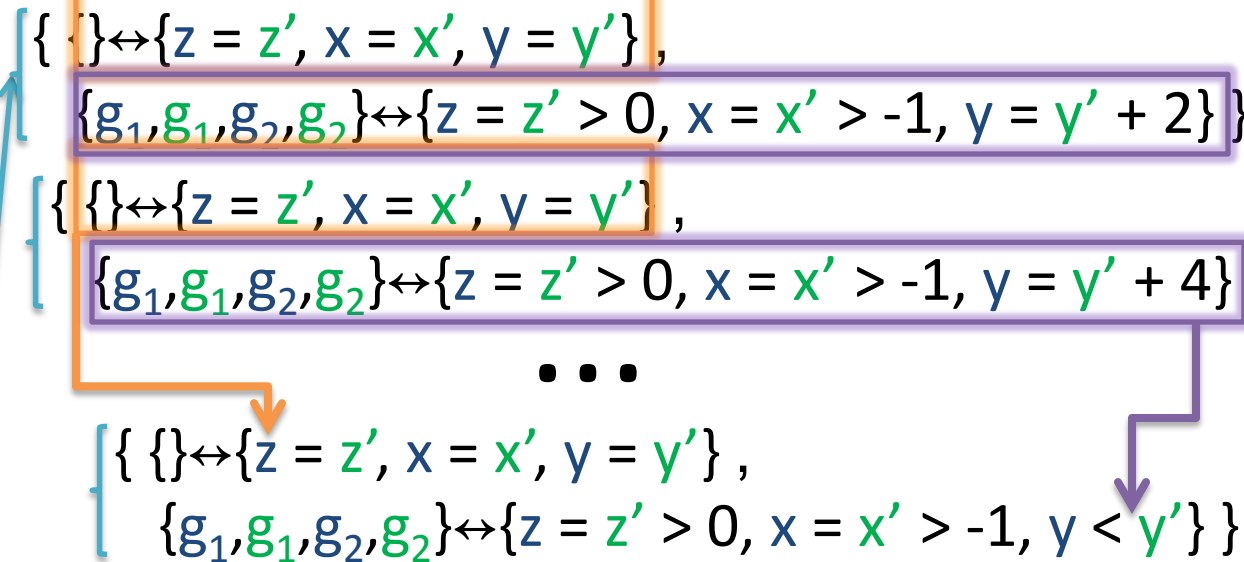
```
    if (g1) x--;
```

```
    if (g'1) x'--;
```

```
    if (g1) goto loop;
```

```
    if (g'1) goto loop';
```

```
}
```



# Results summary

Name	#LOC	#P	Widen	Octagon (Part-by-equiv)	Octagon (Part-by-guard)	Polyhedra (Part-by-equiv)	Polyhedra (Part-by-guard)
remove	16	4	No	✓	✓	✓	✓
copy	44	2	No	✓	✓	✓	✓
fmt	42	5	Yes	✗	TO	✓	✓
md5sum	40	3	Yes	✓	TO	✓	✓
pr	100	10	Yes	TO	TO	✓	TO
savewd	86	1	No	✓	✓	✓	✓
seq	23	15	Yes	✗	TO	✗	✗
addr	77	1	No	✓	TO	✓	TO
nsGDDN	47	11	No	✗	✗	✓	✓
sign	8	2	No	✓	✓	✓	✓
sum	7	5	Yes	✗	✗	✓	✓

# Results

```
bool bsd_split_3 (char *s, size_t s_len,...) {
    int i = s_len;
    i--;

    i = s_len - 1;
    while (i && s[i] != '\0') {
        (1)
        i--;
    }
    (2)
    ...
}
```

$\sigma_1(\text{equivalent})$ :
$s\_len' = s\_len$
$i' = i$
$s\_len' - 1 \geq i'$

$\sigma_1$ :	$\sigma_2(\text{equivalent})$ :
$s\_len = s\_len' = 0$	$s\_len' = s\_len$
$i' \neq i \leq -1$	$i' = i$
	$s\_len' - 1 \geq i'$

```
bool bsd_split_3 (char *s, size_t s_len,...) {
    int i = s_len;
    i--;
    + if (s_len == 0)
    +   return false;
    i = s_len - 1;
    while (i && s[i] != '\0') {
        (1)
        i--;
    }
    (2)
    ...
}
```

# Results

$\sigma 1$	$\sigma 2$	$\sigma 3$
input_position0 = 0	input_position0 < -width	input_position0 < -width
chars' = 0	input_position0 < 0	input_position0 > 0
input_position = width	input_position' = 0	input_position' = 0
input_position < 0	input_position < width	input_position > width
input_position' = 0	width < 0	input_position <= 0

```

int i
bool
int
...

input_position += width;

...
return chars;
}

```

```

...
+ if (width < 0 && input_position == 0) {
+   chars = 0;
+   input_position = 0;
+ } else if (width < 0 && input_position <= -width) {
+   input_position = 0;
+ } else {
+   input_position += width;
+ }
...
return chars;
}

```



# Summary

- Abstract semantic differencing
  - Characterize difference between similar programs
  - Many applications for computed difference
- Key ideas that make this work
  - Correlating program
  - Correlating abstract domain
    - Partially disjunctive based on equivalences
- Results over real-world patches

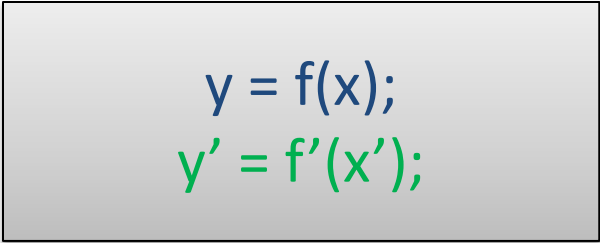
# Proposed Questions

1. Why not use a correlating semantics instead of  $P \bowtie P'$ ?
2. Partitioning
  1. are guard values considered in partitioning?
  2. isn't  $2^{|VAR|}$  too big?
3. How do you handle other data types?
4. How did you handle procedures?
5. How are the variables matched?
6. Future work?

# Backup slides

# Procedures

- Establish equivalence of procedures bottom-up
- Use the equivalence of callees to establish equivalence of callers
  - No recursion
- Similar to Strichman et. al.



```
y = f(x);  
y' = f'(x');
```

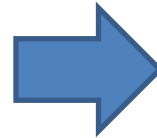
# Related Work

- Brumley et al.
- Strichman et al. (SymDiff)
- DSE
- UC-KLEE

# Sequential composition is bad

```
int sign(int x) {  
  int sgn;  
  if (x < 0)  
    sgn = -1  
  else  
    sgn = 1  
  retval = sgn  
}
```

```
int sign'(int x') {  
  int sgn';  
  if (x' < 0)  
    sgn' = -1  
  else  
    sgn' = 1  
  if (x' == 0)  
    sgn' = 0  
  retval' = sgn'  
}
```



```
int sign; sign'(int x) {  
  int x' = x;  
  int sgn;  
  if (x < 0)  
    sgn = -1  
  else  
    sgn = 1  
  retval = sgn  
  int sgn';  
  if (x' < 0)  
    sgn' = -1  
  else  
    sgn' = 1  
  if (x' == 0)  
    sgn' = 0  
  retval' = sgn'  
  assert(retval == retval')  
}
```

P and P' use disjoint parts of memory

Any interleaving is sufficient (POR)

Can use sequential composition

# Challenge: P;P' Inhibits Partial Disjunction

```
int sign;sign'(int x) {
```

```
  int x' = x;
```

```
  int sgn;
```

```
  if (x < 0)
```

```
    sgn = -1
```

```
  else
```

```
    sgn = 1
```

```
  retval = sgn
```

```
  int sgn';
```

```
  if (x' < 0)
```

```
    sgn' = -1
```

```
  else
```

```
    sgn' = 1
```

```
  if (x' == 0)
```

```
    sgn' = 0
```

```
  retval' = sgn'
```

```
  assert(retval == retval')
```

```
}
```

states that hold equivalence for the same set of variables will be merged

$\sqcup_{\text{eqv}} =$

$\rightarrow [ \{x = x' < 0, \text{sgn} \mapsto -1, \text{sgn}' \mapsto ?\} ]$

$\rightarrow [ \{x = x' \geq 0, \text{sgn} \mapsto 1, \text{sgn}' \mapsto ?\} ]$

$[ \{x = x' \geq 0, \text{sgn} \mapsto [-1,1], \text{sgn}' \mapsto ?\} ]$



**Challenge:** How to compose P, P' such that  $\sqcup_{\text{eqv}}$  will maintain equivalence?

# Correlating Analysis for P;P'

```
int sign;sign'(int x) {
```

```
  int x' = x;
```

```
  int sgn;
```

```
  if (x < 0)
```

```
    sgn = -1 → {x = x' < 0, sgn ↦ -1, sgn' ↦ ?}
```

```
  else
```

```
    sgn = 1 → {x = x' ≥ 0, sgn ↦ 1, sgn' ↦ ?}
```

```
  retval = sgn
```

```
  int sgn';
```

```
  if (x' < 0)
```

```
    sgn' = -1
```

```
  else
```

```
    sgn' = 1
```

```
  if (x' == 0)
```

```
    sgn' = 0
```

```
  retval' = sgn'
```

```
  assert(retval == retval')
```

```
}
```

⊔ =



{x = x', sgn ↦ [-1,1], sgn' ↦ ?}



**Challenge:** Performing a join operation in the correlating domain results in a **non-restorable loss of equivalence**



# Solution: Disjunctive Correlating Domain

```

int sign; sign'(int x) {
  int x' = x;
  int sgn;
  if (x < 0)
    sgn = -1 → [ {x = x' < 0, sgn ↦ -1, sgn' ↦ ?} ]
  else
    sgn = 1 → [ {x = x' ≥ 0, sgn ↦ 1, sgn' ↦ ?} ]
  retval = sgn
  int sgn';
  if (x' < 0)
    sgn' = -1
  else
    sgn' = 1
  if (x' == 0)
    sgn' = 0
  retval' = sgn'
  assert(retval == retval') → [ {x = x' < 0, sgn = sgn' ↦ -1}, {x = x' > 0, sgn = sgn' ↦ 1},
                                {x = x' ↦ 0, sgn ↦ 1, sgn' ↦ 0} ]
}

```

**U =**  
 ↓  
 [ {x = x' < 0, sgn ↦ -1, sgn' ↦ ?}, {x = x' ≥ 0, sgn ↦ 1, sgn' ↦ ?} ]  
 . . .

**Challenge:** Every condition splits the abstract state, does not scale

# Patch

- (Why) Easier to spot the difference
  - Our approach's precision is based on syntactic similarity
- (What) A patch can be any syntactic change
  - can be fragments or compound patch
  - the smaller the better

