

A Theory For Semantic Type Checking

Motivation: if we compile/optimize a typed program we are often unable to use the source-level type system to type check the target program.

Examples:

① after register allocation a given variable is used with different types

```
{(x: int) }
```

```
  x ← x == int 7;
```

```
  x ← not x
```

```
  assert (x: bool) ← postcondition
```

```
}
```

• we want to check (just like in the source) that operations are used with ~~the~~ appropriate operands:

== int	both operands are int and result is bool
not	operand and result is bool

- we cannot use a type system based on declarations
- we need to use a data-flow based type checker
- or VGen!

Assume that we have predicates of the form
 $e : \tau$ (for now $\tau \in \{\text{int}, \text{bool}\}$)

$\text{fun}(x : \text{int})$ $x \leftarrow x == \text{int} \quad \nabla$ $x \leftarrow \text{not } x$ $\text{assert}(x : \text{bool})$	$\text{VC} = x : \text{int} \wedge x == \nabla : \text{bool} \wedge \text{not}(x == \nabla) : \text{bool}$ $\text{VC} = x : \text{bool} \wedge \text{not } x : \text{bool}$ $\text{VC} = x : \text{bool}$
---	---

The VC for the whole function is.

$$\forall x. x : \text{int} \Rightarrow x : \text{int} \wedge x == \nabla : \text{bool} \wedge \text{not}(x == \nabla) : \text{bool}$$

which we can prove using some simple inference rule that model the source-level type rules

$$\frac{E_1 : \text{int} \quad E_2 : \text{int}}{E_1 ==_{\text{int}} E_2 : \text{bool}}$$

$$\frac{E : \text{bool}}{\text{not } E : \text{bool}}$$

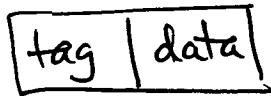
Example 2

In the process of compilation some source-level abstraction is eliminated, exposing concrete representation details

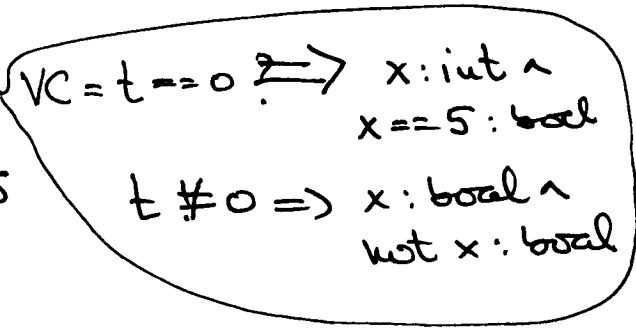
- consider union types:
- concrete representation

```
p: int + bool
case p of
  left(x) → x == 5
  | right(x) → not x
```

p is a pointer to



```
t ← *p
x ← *(p+1)
if t == 0 then
  res ← x == 5
else
  res ← not x
assert (res : bool)
```



- we cannot assign a source-level type to x before the conditional.
- but a data-flow based approach can keep track of conditionals
- VCGen can do that too.

We will want to maintain a few invariants such as

- only 0 and 1 are valid representations for booleans
- if we have a τ^* then we can read and write (a τ) to that address.

$$\frac{}{0: \text{bool}} \quad \frac{}{1: \text{bool}} \quad \frac{e: \tau \quad e=e'}{e': \tau}$$

We try the rule

$$\frac{e: \tau^*}{\text{sel}(m, e): \tau}$$

but this is not sound!!

It would allow us to prove

$$\text{sel}(\text{upd}(m, e, 0), e): \tau \quad \text{from } e: \tau^*$$

from which we can prove $0: \tau$!!

A big hole in the type system

The congruence rule gives us power but it destroys any abstraction!! Because of that we have to be very explicit about assumptions

We have to add an assumption that the contents of memory matches the types.

$\text{typed}(m)$

We need first to define the semantics of typing predicates

• we define a representation function R that associates with each base type τ a set of values in U that are valid representations of expressions of type τ

U is the universe of all base values (machine words in assembly language)

Not all R are good. We define a set of conditions that each R must ~~hold~~ satisfy.

a) $R(\text{int}) = U$ ← all values are valid represent. of integers.

b) $R(\text{bool}) = \{0, 1\}$

c) ~~$v \in R(\sigma^*) \Rightarrow$~~
 $v \in R((\sigma_1 \times \sigma_2)^*) \Rightarrow v \in R(\sigma_1^*)$ and
 ~~$v \in R(\sigma_2^*)$~~
 $v + |\sigma_1| \in R(\sigma_2^*)$

We define the size of a structured type

$$|\tau| = 1 \quad |\sigma_1 \times \sigma_2| = |\sigma_1| + |\sigma_2|$$

$$|\sigma[]| = 2 \quad |\text{int. } \sigma| = |\sigma| \quad |\sigma_1 + \sigma_2| = 1 + \max(|\sigma_1|, |\sigma_2|)$$

• is not defined for type variables. All t must occur inside a σ^* or $\sigma[l]$

$$d) \mathcal{R}((\sigma_1 + \sigma_2)^*) = \mathcal{R}((\text{const}(0) \times \sigma_1)^*) \cup \mathcal{R}((\text{const}(1) \times \sigma_2)^*)$$

$$e) \mathcal{R}(\text{const}(n)) = \{n\}$$

$$f) v \in \mathcal{R}(\sigma[l]) \text{ then } \forall i: 0 \leq i < l \\ v + i \cdot |\sigma| \in \mathcal{R}(\sigma^*)$$

$$g) \mathcal{R}(\text{[ut.}\sigma\text{]}^*) = \mathcal{R}(\text{[ut.}\sigma/t\text{]}\sigma^*)$$

$$h) v \in \mathcal{R}(\sigma[l]^*) \text{ then } v \notin \mathcal{R}(\tau^*) \text{ for any } \tau$$

$$i) v \in \mathcal{R}(\tau^*) \text{ and } v \in \mathcal{R}(\tau'^*) \text{ then } \tau \equiv \tau'$$

↑ this classifies all memory locations uniquely

And the meaning of $\text{typed}(m)$ in a given representation

$$a) v \in \mathcal{R}(\tau^*) \text{ then } v \in \text{Dom}(m) \text{ and } m(v) \in \mathcal{R}(\tau)$$

$$b) v \in \mathcal{R}(\sigma[l]^*) \text{ then } v, v+1 \in \text{Dom}(m) \text{ and}$$

$$m(v) > 0 \text{ and}$$

$$m(v+1) \in \mathcal{R}(\sigma[m(v)])$$

Now we can give the inference rules

$$\frac{A : \sigma[L] \quad I \geq 0 \quad I < L \quad |\sigma| = S}{S * I + A : \sigma^*}$$

$$\frac{A : (\sigma_1 \times \sigma_2)^* \quad |\sigma_1| = S_1}{S_1 + A : \sigma_2^*} \quad \#$$

$$\frac{A : (\sigma_1 \times \sigma_2)^*}{A : \sigma_1^*}$$

$$\frac{A : \tau^* \quad \text{typed}(M)}{\text{sel}(M, A) : \tau}$$

(works for a τ
but not for a σ)

~~$$\frac{\text{[scribble]}}{E : \text{const}(E^*)}$$~~

$$\frac{A : (\sigma[])^* \quad \text{typed}(M)}{\text{sel}(M, A+1) : \sigma[\text{sel}(M, A)]}$$

$$\frac{}{0 : \text{bool}}$$

$$\frac{}{1 : \text{bool}}$$

$$\frac{A : (\sigma_1 + \sigma_2) * \quad \text{typed}(M) \quad \text{sel}(M, A) = 0}{A+1 : \sigma_1 *}$$

$$\frac{\text{typed}(M) \quad A : \tau * \quad E : \tau}{\text{typed}(\text{upd}(M, A, E))} \quad \text{upd}$$

We can even prove the soundness of these rules.
Take the upd rule

$$\text{Hyp: } R \Vdash \text{typed}(M) \quad A \in \mathcal{R}(\tau *) \quad E \in \mathcal{R}(\tau)$$

Need to prove

$$a) \forall v \in \mathcal{R}(\tau' *) \quad v \in \text{Dom}(m') \quad m'(v) \in \mathcal{R}(\tau')$$

$$\text{where } m' = \text{upd}(m, A, E)$$

$$a.1. \quad v = A \quad \text{then, by inv. i} \quad \tau = \tau'$$

$$m'(v) = E \quad \checkmark$$

$$a.2. \quad v \neq A \quad \text{then } m'(v) = m(v) \quad \checkmark$$

$$(\text{Dom}(m') = \text{Dom}(m))$$

$$b) \forall v \in \mathcal{R}((\sigma[\])*) \quad \text{then } v, v+1 \in \text{Dom}(m')$$

$$m'(v) > 0$$

$$m'(v+1) \in \mathcal{R}((\sigma[m(v)])*)$$

$$\text{but } v \neq E \quad \Rightarrow \quad m'(v) = m(v) > 0$$

Now $v+1 \in \mathcal{R}((\sigma[m(v)])*)$ and by a) above
we have q.e.d.

A proof procedure for typing predicates

- tactic based, reversing the inference rules
- cannot just do backwards chaining because there are vars in hyp. that do not appear in conclusion.

· We need a function that given A it finds out all potential σ such that $A : \sigma^*$

The idea is that pointers can be created in a limited nr. of ways

A is $A' + S$ · S is a constant

- second element of a pair
- index in an array

A is $\text{sel}(M, A')$ · try to prove that $A' : \tau^*$

A is variable · we must find the type of A within assumptions.

· This is not going to be complete (what if a pointer is computed in a convoluted way

```
{ (int **)  
  * (x + alwaysZero());  
}
```