

Typed Assembly Language

Evvari Koppel and Magnus Leppik

About presentation

- Why type checking?
- Proof-carrying code
- Typed Assembly Language
- TAL-0
- TAL-1
- Examples

What is TAL

- TAL is assembly language that is extended to make use of annotating datatypes for each used value
- Type checker – to check how code acts when executed
- Type safety – prevent type errors

Why Typing Checking

- Low-level code and high-level program
- Type checking – convenient way to ensure that a program has certain semantic properties
- Major component of the security infrastructure in distributed systems
- Memory safety
- Type safety
- Malicious code

Proof-carrying code (1)

- The principle of PCC is that the need to trust a piece of code is eliminated by machine-checkable proof that the code has certain properties.
- Using PCC to build trustworthy systems:
 - What properties should we require of the code?
 - How do code producers construct a formal proof that their code has the desired properties?

Proof-carrying code (2)

- Solution: type-preserving compilation.
- We seek a principled approach to the design of typed intermediate language.

TAL-0: Control-Flow-Safety(1)

- Control-Flow safety
- Focus on control-flow safety will let us start with simple abstract machine
- The syntax for control-flow-safety assembly language:

$r ::=$ r1 r2 ... rk	<i>registers:</i>	$t ::=$	<i>instructions:</i>
$v ::=$ n ℓ r	<i>operands:</i> integer literal label or pointer registers	$r_d := v$ $r_d := r_s + v$ if r jump v	
		$I ::=$ jump v $t; I$	<i>instruction sequences:</i>

TAL-0: Control-Flow-Safety(2)

- We model evaluation of TAL-0 assembly programs using a rewriting relation between *abstract* machine states.
- We maintain the distinction between labels and arbitrary integers.
- Enforcing the safety property now reduces to ensuring that abstract machine cannot get stuck.

TAL-0: Control-Flow-Safety(3)

- Syntax for TAL-0 abstract machines:

$\mathcal{R} ::=$ $\{r_1 = v_1, \dots, r_k = v_k\}$	<i>register files:</i>	$H ::=$	<i>heaps:</i>
$\hat{h} ::=$	<i>heap values:</i>	$\{\ell_1 = h_1, \dots, \ell_m = h_m\}$	
I	<i>code</i>	$M ::=$	<i>machine states:</i>
		(H, R, I)	

TAL-0: Control-Flow-Safety(4)

- Rewriting rules for TAL-0:

$$\frac{H(\hat{R}(v)) = I}{(H, R, \text{jump } v) \rightarrow (H, R, I)} \quad \text{(JUMP)}$$
$$(H, R, r_d := v; I) \rightarrow (H, R[r_d = \hat{R}(v)], I) \quad \text{(MOV)}$$
$$\frac{R(r_s) = n_1 \quad \hat{R}(v) = n_2}{(H, R, r_d := r_s + v; I) \rightarrow (H, R[r_d = n_1 + n_2], I)} \quad \text{(ADD)}$$
$$\frac{R(r) = 0 \quad H(\hat{R}(v)) = I'}{(H, R, \text{if } r \text{ jump } v; I) \rightarrow (H, R, I')} \quad \text{(IF-EQ)}$$
$$\frac{R(r) = n \quad n \neq 0}{(H, R, \text{if } r \text{ jump } v; I) \rightarrow (H, R, I)} \quad \text{(IF-NEQ)}$$

TAL-0 Type System (1)

- Goal: ensure that any well-formed abstract machine M cannot get stuck.
- Our type system has to:
 - Distinguish labels from integers
 - Ensures that operands of a control transfer are labels
 - No matter how many steps are taken by M, it never gets into a stuck state (i.e typing preserved)

TAL-0 Type System (2)

- Type syntax:

$\tau ::=$	<i>operand types:</i>	$\Gamma ::=$	<i>register file types:</i>
int	<i>word-sized integers</i>	$\{r_1 : \tau_1, \dots, r_k : \tau_k\}$	
code(Γ)	<i>code labels</i>	$\Psi ::=$	<i>heap types:</i>
α	<i>type variables</i>	$\{\ell_1 : \tau_1, \dots, \ell_n : \tau_n\}$	
$\forall \alpha. \tau$	<i>universal polymorphic types</i>		

TAL-0 Type System (3)

- We now formalize the type system using the inference rules:

<p><i>Values</i></p> $\frac{}{\Psi \vdash n : \text{int}} \quad (\text{S-INT})$ $\frac{}{\Psi \vdash \ell : \Psi(\ell)} \quad (\text{S-LAB})$ <p><i>Operands</i></p> $\frac{}{\Psi; \Gamma \vdash r : \Gamma(r)} \quad (\text{S-REG})$ $\frac{\Psi \vdash v : \tau}{\Psi; \Gamma \vdash v : \tau} \quad (\text{S-VAL})$ $\frac{\Psi; \Gamma \vdash v : \forall \alpha. \tau}{\Psi; \Gamma \vdash v : \tau[\tau'/\alpha]} \quad (\text{S-INST})$ <p><i>Instructions</i></p> $\frac{\Psi; \Gamma \vdash v : \tau}{r_d := v : \Gamma \rightarrow \Gamma[r_d : \tau]} \quad (\text{S-MOV})$ $\frac{r_s : \text{int} \quad \Psi; \Gamma \vdash v : \text{int}}{r := r_s + v : \Gamma \rightarrow \Gamma[r_d : \text{int}]} \quad (\text{S-ADD})$ $\frac{\text{int} \quad \Psi; \Gamma \vdash v : \text{code}(\Gamma)}{\text{if } r_s \text{ jump } v : \Gamma \rightarrow \Gamma} \quad (\text{S-IF})$	$\Psi \vdash v : \tau$	<p><i>Instruction Sequences</i></p> $\frac{\Psi; \Gamma \vdash v : \text{code}(\Gamma)}{\Psi \vdash \text{jump } v : \text{code}(\Gamma)} \quad (\text{S-JUMP})$ $\frac{\Psi \vdash \iota : \Gamma \rightarrow \Gamma_2 \quad \Psi \vdash I : \text{code}(\Gamma_2)}{\Psi \vdash \iota; I : \text{code}(\Gamma)} \quad (\text{S-SEQ})$ $\frac{\Psi \vdash I : \tau}{\Psi \vdash I : \forall \alpha. \tau} \quad (\text{S-GEN})$ <p><i>Register Files</i></p> $\frac{\forall r. \Psi \vdash R(r) : \Gamma(r)}{\Psi \vdash R : \Gamma} \quad (\text{S-REGFILE})$ <p><i>Heaps</i></p> $\frac{\forall \ell \in \text{dom}(\Psi). \Psi \vdash H(\ell) : \Psi(\ell) \quad \text{FTV}(\Psi(\ell)) = \emptyset}{\vdash H : \Psi} \quad (\text{S-HEAP})$ <p><i>Machine States</i></p> $\frac{\vdash H : \Psi \quad \Psi \vdash R : \Gamma \quad \Psi \vdash I : \text{code}(\Gamma)}{\vdash (H, R, I)} \quad (\text{S-MACH})$	$\Psi \vdash I : \tau$
---	------------------------	--	------------------------

Proof of Type Soundness for TAL-0 (1)

- It suffices to show:
 - Well-typed machine state is not immediately stuck (progress)
 - When it steps to a new machine state M' , that state is also well-typed (preservation).

Proof of Type Soundness for TAL-0 (2)

LEMMA [TYPE SUBSTITUTION]: If:

1. $\Psi; \Gamma \vdash v : \tau_1$, then $\Psi; \Gamma[\tau/\alpha] \vdash v : \tau_1[\tau/\alpha]$.
2. $\Psi \vdash \iota : \Gamma_1 \rightarrow \Gamma_2$ then $\Psi \vdash \iota : \Gamma_1[\tau/\alpha] \rightarrow \Gamma_2[\tau/\alpha]$.
3. $\Psi \vdash I : \tau_1$, then $\Psi \vdash I : \tau_1[\tau/\alpha]$.
4. $\Psi \vdash R : \Gamma$, then $\Psi \vdash R : \Gamma[\tau/\alpha]$.

Proof of Type Soundness for TAL-0 (3)

LEMMA [REGISTER SUBSTITUTION]: If $\vdash H : \Psi$, $\Psi \vdash R : \Gamma$ and $\Psi; \Gamma \vdash v : \tau$ then $\Psi; \Gamma \vdash \hat{R}(v) : \tau$ \square

LEMMA [CANONICAL VALUES]: If $\vdash H : \Psi$ and $\Psi \vdash v : \tau$ then:

1. If $\tau = \text{int}$ then $v = n$ for some n .
2. If $\tau = \text{code}(\Gamma)$ then $v = \ell$ for some $\ell \in \text{dom}(H)$ and $\Psi \vdash H(\ell) : \text{code}(\Gamma)$. \square

Proof of Type Soundness for TAL-0 (4)

LEMMA [CANONICAL OPERANDS]: If $\vdash H : \Psi$, $\Psi \vdash R : \Gamma$, and $\Psi; \Gamma \vdash v : \tau$ then:

1. If $\tau = \text{int}$ then $\hat{R}(v) = n$ for some n .
2. If $\tau = \text{code}(\Gamma)$ then $\hat{R}(v) = \ell$ for some $\ell \in \text{dom}(H)$ and $\Psi \vdash H(\ell) : \text{code}(\Gamma)$. \square

THEOREM [SOUNDNESS OF TAL-0]: If $\vdash M$, then there exists an M' such that $M \rightarrow M'$ and $\vdash M'$. \square

Proof Representation and Checking

- For TAL-0 it is sufficient to provide types for the labels;
- Keep the type checker as simple as possible:
 - a) Type reconstruction is entirely syntax directed (for any given term at most one rule should apply)
 - b) Explicit representation of the complete proof of well-formedness
 - We can ship the binary machine code, disassemble it and then compare it against the assembly-level proof (proof-carrying code)

TAL-1: Simple Memory-Safety (1)

- TAL-0 includes registers and heap-allocated code; no support for allocated *data*.
- TAL-1:
 - adds primitive support for allocated objects that can be shared by reference (i.e pointer)
 - includes a notion of object-level memory safety.
- How to accomodate locations that hold values of different types at different times?

TAL-1: Simple Memory-Safety (2)

```
{r1:ptr(code(...))}  
1.  r3 := 0;  
2.  Mem[r1] := r3;  
3.  r4 := Mem[r1];  
4.  jump r4
```

- The code above should be rejected by the type-checker (control-flow safety property)

```
{r1:ptr(code(...)),r2:ptr(code(...))}  
1.  r3 := 0;  
2.  Mem[r1] := r3;  
3.  r4 := Mem[r2];  
4.  jump r4
```

TAL-1: Simple Memory-Safety (3)

- We need some support for
 - Allocating and initializing data structures that are to be shared;
 - Stack-allocating procedure frames.
- Separate locations into two classes:
 - Shared pointers that support arbitrary aliasing;
 - Unique pointers that will support updates that change the type of the contents.

The TAL-1 Extended Abstract Machine (1)

Syntactic extensions to TAL-0 and rewriting rules:

$r ::=$	<i>registers:</i>
$r_1 \mid r_2 \mid \dots \mid r_k$	<i>gp registers</i>
	<i>stack pointer</i>
$i ::=$	<i>instructions:</i>
...	<i>as in TAL-0</i>
$r := \text{Mem}[r_s + n]$	<i>load from memory</i>
$\text{Mem}[r_d + n] := r_s$	<i>store to memory</i>
$r := \text{malloc } n$	<i>allocate n heap words</i>
$\text{Commit } r_d$	<i>become shared</i>
$\text{salloc } n$	<i>allocate n stack words</i>

The TAL-1 Extended Abstract Machine (2)

	$sfree\ n$	<i>free n stack words</i>
$v ::=$		<i>operands:</i>
	r	<i>registers</i>
	n	<i>integer literals</i>
	ℓ	<i>code or shared data pointers</i>
	$uptr(h)$	<i>unique data pointers</i>
$h ::=$		<i>heap values:</i>
	I	<i>instruction sequences</i>
	$\langle v_1, \dots, v_n \rangle$	<i>tuples</i>

The TAL-1 Extended Abstract Machine (3)

- The rewriting rules for the instructions of TAL-1

$$\frac{\hat{R}(v) \neq \text{uptr}(h)}{(H, R, r_d := v; I) \rightarrow (H, R[r_d = v], I)} \quad (\text{MOV-1})$$

This rule can only fire when the source operand is not a unique pointer.

We must now give the rewriting rules for the new instructions:

$$(H, R, r_d := \text{malloc } n; I) \rightarrow (H, R[r_d = \text{uptr}\langle m_1, \dots, m_n \rangle], I) \quad (\text{MALLOC})$$

$$\frac{r_d \neq \text{sp} \quad \ell \notin \text{dom}(H)}{(H, R[r_d = \text{uptr}(h)], \text{commit } r_d; I) \rightarrow (H[\ell = h], R[r_d = \ell], I)} \quad (\text{COMMIT})$$

The TAL-1 Extended Abstract Machine (4)

$$\frac{R(r_s) = \ell \quad H(\ell) = \langle v_0, \dots, v_n, \dots, v_{n+m} \rangle}{(H, R, r_d := \text{Mem}[r_s + n]; I) \rightarrow (H, R[r_d = v_n], I)} \quad (\text{LD-S})$$

$$\frac{R(r_s) = \text{uptr} \langle v_0, \dots, v_n, \dots, v_{n+m} \rangle}{(H, R, r_d := \text{Mem}[r_s + n]; I) \rightarrow (H, R[r_d = v_n], I)} \quad (\text{LD-U})$$

$$\frac{R(r_d) = \ell \quad H(\ell) = \langle v_0, \dots, v_n, \dots, v_{n+m} \rangle \quad R(r_s) = v \quad v \neq \text{uptr}(h)}{(H, R, \text{Mem}[r_d + n] := r_s; I) \rightarrow (H[\ell = \langle v_0, \dots, v, \dots, v_{n+m} \rangle], R, I)} \quad (\text{ST-S})$$

The TAL-1 Extended Abstract Machine (5)

$$\frac{R(r_d) = \text{uptr}\langle v_0, \dots, v_n, \dots, v_{n+m} \rangle, \quad R(r_s) = v \quad v \neq \text{uptr}\langle n \rangle}{(H, R, \text{Mem}[r_d + n] := r_s; I) \rightarrow (H, R[r_d = \text{uptr}\langle v_0, \dots, v, \dots, v_{n+m} \rangle], I)} \quad (\text{ST-U})$$

$$\frac{R(\text{sp}) = \text{uptr}\langle v_0, \dots, v_p \rangle \quad p + n \leq \text{MAXSTACK}}{(H, R, \text{salloc } n) \rightarrow (H, R[\text{sp} = \text{uptr}\langle m_1, \dots, m_n, v_0, \dots, v_p \rangle])} \quad (\text{SALLOC})$$

$$\frac{R(\text{sp}) = \text{uptr}\langle m_1, \dots, m_n, v_0, \dots, v_p \rangle}{(H, R, \text{sfree } n) \rightarrow (H, R[\text{sp} = \text{uptr}\langle v_0, \dots, v_p \rangle])} \quad (\text{SFREE})$$

TAL-1 Changes to the Type System (1)

- New set of types for classifying TAL-1 values and new typing rules:

$\tau ::=$	<i>operand types:</i>	$\sigma ::=$	<i>allocated types:</i>
...	<i>as in TAL-0</i>	ϵ	<i>empty</i>
$\text{ptr}(\sigma)$	<i>shared data pointers</i>	τ	<i>value type</i>
$\text{uptr}(\sigma)$	<i>unique data pointers</i>	σ_1, σ_2	<i>adjacent</i>
$\forall \rho, \tau$	<i>quantification over allocated types</i>	ρ	<i>allocated type variable</i>

TAL-1 Changes to the Type System (2)

- New typing rules:

Heap Values	$\Psi \vdash v : \tau$	Operands	$\Psi \vdash v : \tau$
$\Psi; \Gamma \vdash v_i : \tau_i$		$\Psi; \Gamma \vdash h : \sigma$	
<hr/>		<hr/>	
$\Psi \langle v_1, \dots, v_n \rangle : \tau_1, \dots, \tau_n$	(S-TUPLE)	$\Psi; \Gamma \vdash \text{uptr}(h) : \text{uptr}(\sigma)$	(S-UPTR)

TAL-1 Changes to the Type System (3)

- New typing rules:

Instructions

$\Psi \vdash t : \Gamma_1 \rightarrow \Gamma_2$

$$\frac{\Psi; \Gamma \vdash v : \tau \quad \tau \neq \text{uptr}(\sigma)}{\Psi \vdash r_d := v : \Gamma \rightarrow \Gamma[r_d : \tau]}$$

(S-MOV-1)

$$n \geq 0$$

$$\Psi \vdash r_d := \text{malloc } n : \Gamma \rightarrow \Gamma[r_d : \text{uptr}(\underbrace{\text{int}, \dots, \text{int}}_n)]$$

(S-MALLOC)

$$\frac{\Psi; \Gamma \vdash r_d : \text{uptr}(\sigma) \quad r_d \neq \text{sp}}{\Psi \vdash \text{commit } r_d : \Gamma \rightarrow \Gamma[r_d : \text{ptr}(\sigma)]}$$

(S-COMMIT)

$$\frac{\Psi; \Gamma \vdash r_s : \text{ptr}(\tau_1, \dots, \tau_n, \sigma)}{\Psi \vdash r_d := \text{Mem}[r_s + n] : \Gamma \rightarrow \Gamma[r_d : \tau_n]}$$

(S-LDS)

$$\frac{\Psi; \Gamma \vdash r_s : \text{uptr}(\tau_1, \dots, \tau_n, \sigma)}{\Psi \vdash r_d := \text{Mem}[r_s + n] : \Gamma \rightarrow \Gamma[r_d : \tau_n]}$$

(S-LDU)

TAL-1 Changes to the Type System (4)

$$\begin{array}{c}
 \frac{\Psi; \Gamma \vdash r_d \geq \tau_d \quad \tau_d \neq \text{uptr}(\sigma) \quad \Psi; \Gamma \vdash r_d : \text{ptr}(\tau_1, \dots, \tau_n, \sigma)}{\Psi \vdash \text{Mem}[r_d + n] := r_d : \Gamma \rightarrow \Gamma} \quad (\text{S-STG}) \\
 \\
 \frac{\Psi; \Gamma \vdash r_d : \tau \quad \tau \neq \text{uptr}(\sigma) \quad \Psi; \Gamma \vdash r_d : \text{uptr}(\tau_1, \dots, \tau_n, \sigma)}{\Psi \vdash \text{Mem}[r_d + n] := r_d : \Gamma \rightarrow \Gamma[r_d : \text{uptr}(\tau_1, \dots, \tau_n, \sigma)]} \quad (\text{S-STU}) \\
 \\
 \frac{\Psi; \Gamma \vdash sp : \text{uptr}(\sigma) \quad n \geq 0}{\Psi \vdash \text{salloc } n : \Gamma \rightarrow \Gamma[\underbrace{sp : \text{uptr}(\text{int}, \dots, \text{int}, \sigma)}_n]} \quad (\text{S-SALLOC}) \\
 \\
 \frac{\Psi; \Gamma \vdash sp : \text{uptr}(\tau_1, \dots, \tau_n, \sigma)}{\Psi \vdash \text{sfree } n : \Gamma \rightarrow \Gamma[\text{sp : uptr}(\sigma)]} \quad (\text{S-SFREE})
 \end{array}$$

- At this point TAL-1 provides enough mechanism for the compiler of a polymorphic, procedural language.

Compiling to TAL-1 (1)

- A simple example:

```
int prod (int x, int y){  
  int a = 0;  
  while (x != 0) {  
    a = a + y;  
    x = x - 1;  
  }  
  return a;  
}
```

Compiling to TAL-1 (2)

```
prod:  $\forall a,b,c,s.$   
  code{r1:a,r2:b,r3:c,sp:uptr(int,int,s),  
    r4: $\forall d,e,f.$ code{r1:int,r2:d,r3:e,r4:f,sp:uptr(s)}}  
  r2 := Mem[sp]; // r2:int, r2 := x  
  r3 := Mem[sp+1]; // r3:int, r3 := y  
  r1 := 0 // r1:int, a := 0  
  jump loop  
  
loop:  $\forall s.$ code{r1,r2,r3:int,sp:uptr(int,int,s),  
  r4: $\forall d,e,f.$ code{r1:int,r2:d,r3:e,r4:f,sp:uptr(s)}}  
  if r2 jump done; // if x  $\leftrightarrow$  0 goto done  
  r1 := r1 + r3; // a := a + y  
  r2 := r2 + (-1); // x := x - 1  
  jump loop  
  
done:  $\forall s.$ code{r1,r2,r3:int,sp:uptr(int,int,s),  
  r4: $\forall d,e,f.$ code{r1:int,r2:d,r3:e,r4:f,sp:uptr(s)}}  
  sfree 2; // sp:uptr(s)  
  jump r4
```

Compiling to TAL-1 (3)

```
prod:  $\forall a,b,c,s.$   
  code{r1:a,r2:b,r3:c,sp:uptr(int,int,s),  
    r4: $\forall d,e,f.$ code{r1:int,r2:d,r3:e,r4:f,sp:uptr(s)}}  
r2 := Mem[sp];    // r2:int, r2 := x  
r3 := Mem[sp+1];  // r3:int, r3 := y  
r1 := 0           // r1:int, a := 0  
jump loop
```

Compiling to TAL-1 (4)

```
loop:  $\forall s.$ code{r1,r2,r3:int,sp:uptr(int,int,s),  
           r4: $\forall d,e,f.$ code{r1:int,r2:d,r3:e,r4:f,sp:uptr(s)}}  
if r2 jump done; // if x  $\leftrightarrow$  0 goto done  
r1 := r1 + r3;    // a := a + y  
r2 := r2 + (-1); // x := x - 1  
jump loop
```

Compiling to TAL-1 (5)

```
done:  $\forall s.$ code{r1,r2,r3:int,sp:uptr(int,int,s),  
          r4: $\forall d,e,f.$ code{r1:int,r2:d,r3:e,r4:f,sp:uptr(s)}}  
sfree 2;      // sp:uptr(s)  
jump r4
```

Some Real World Issues

- TAL-1 and the extensions described earlier provide mechanisms needed to implement only very simple languages.
- Further extensions to TAL-1:
 - STAL;
 - TALT;
 - TALx86.

Conclusion

- The typing annotations are produced and consumed by machines;
- Low-level languages present new challenges to type system designers;
- Ideally, proofs should be carried out in a machine-checked environment.

Used materials

Benjamin C. Pierce edition

"Advanced Topics in Types and Programming Languages " (MIT Press, 2005)

– Typed Assembly Language, by Greg Morrisett

Thank you!