# PROOF-CARRYING-CODE

### APPLYING FORMAL METHODS IN A DISTRIBUTED WORLD

Hans-Wolfgang Loidl

LFE Theoretische Informatik, Institut für Informatik,
Ludwig-Maximilians Universität, München

June 30, 2005

**❶ ENCODING PROOFS**

**❷ PROGRAM LOGICS**

**❸ TCB SIZE**

**❹ PCC IN ACTION: CCURED**

## MEETING THE CHALLENGES

The previous lecture explained the concepts behind PCC, its strengths and weaknesses:

- 🙂 Unforgable certificates
- 🙂 Separation of code safety and trust
- 🙁 High overhead in terms of certificate size and/or trusted code base (TCB)

In this lecture we will look into the details of making the components work.

## LF TERMS

The Logical Framework (LF) is a generic description of logics.
Entities on three levels: objects, families of types, and kinds.

$$
\begin{array}{llll}
\textit{Kinds} & K & ::= & \text{Type} \mid \Pi x : A.K \\
\textit{Families} & A & ::= & a \mid \Pi x : A.B \mid \lambda x : A.B \mid A\,M \\
\textit{Objects} & M & ::= & c \mid x \mid \lambda x : A.M \mid M\,N
\end{array}
$$

Signatures: mappings of constants to types and kinds
Contexts: mappings of variables to types

$$
\begin{array}{llll}
\textit{Signatures} & \Sigma & ::= & \langle\rangle \mid \Sigma, a : K \mid \Sigma, c : A \\
\textit{Contexts} & \Gamma & ::= & \langle\rangle \mid \Gamma, x : A
\end{array}
$$

## LF Type System

Judgements:

$$\Gamma \vdash_\Sigma A : K$$

meaning $A$ has kind $K$ in context $\Gamma$ and signature $\Sigma$.

$$\Gamma \vdash_\Sigma M : A$$

meaning $M$ has type $A$ in context $\Gamma$ and signature $\Sigma$.

# LF Type System (objects)

$$\frac{\vdash_\Sigma \Gamma \qquad c : A \in \Sigma}{\Gamma \vdash_\Sigma c : A} \qquad \text{(const-obj)}$$

$$\frac{\vdash_\Sigma \Gamma \qquad x : A \in \Gamma}{\Gamma \vdash_\Sigma x : A} \qquad \text{(var-obj)}$$

$$\frac{\Gamma, x : A \vdash_\Sigma M : B}{\Gamma \vdash_\Sigma \lambda x : A.M : \Pi x : A.B} \qquad \text{(abs-obj)}$$

$$\frac{\Gamma \vdash_\Sigma M : \Pi x : A.B \qquad \Gamma \vdash_\Sigma N : A}{\Gamma \vdash_\Sigma M \; N : [N/x]B} \qquad \text{(app-obj)}$$

$$\frac{\Gamma \vdash_\Sigma M : A \qquad \Gamma \vdash_\Sigma A' : \text{type} \qquad \Gamma \vdash_\Sigma A \equiv A'}{\Gamma \vdash_\Sigma M : A'} \qquad \text{(conv-obj)}$$

## ENCODING THE LOGIC INTO LF

3 LF-level types are used: exp for expressions, pred for predicates, and tp for types.

Encoding constants and terms:

$$+ \ : \ \text{exp} \rightarrow \text{exp} \rightarrow \text{exp}$$
$$\text{true} \ : \ \text{pred}$$
$$\text{impl} \ : \ \text{pred} \rightarrow \text{pred} \rightarrow \text{pred}$$
$$\text{all} \ : \ (\text{pred} \rightarrow \text{pred}) \rightarrow \text{pred}$$
$$\cdots$$

## Encoding the Logic into LF

3 LF-level types are used: exp for expressions, pred for predicates, and tp for types.

Encoding constants and terms:

$$+ \ : \ \mathrm{exp} \rightarrow \mathrm{exp} \rightarrow \mathrm{exp}$$
$$\mathrm{true} \ : \ \mathrm{pred}$$
$$\mathrm{impl} \ : \ \mathrm{pred} \rightarrow \mathrm{pred} \rightarrow \mathrm{pred}$$
$$\mathrm{all} \ : \ (\mathrm{pred} \rightarrow \mathrm{pred}) \rightarrow \mathrm{pred}$$
$$\cdots$$

Note that all is higher order. We can use the application of LF-level types to encode substitution.

## ENCODING THE LOGIC INTO LF

Encoding proofs: $\text{pf} : \text{pred} \rightarrow \text{Type}$

$$\begin{aligned}
\text{and\_i} : \;\; &\Pi \; p\colon \text{pred}. \; \Pi \; r\colon \text{pred}. \\
&\text{pf } p \;\rightarrow\; \text{pf } r \;\rightarrow\; \text{pf } (\text{and } p \; r) \\
\text{all\_i} : \;\; &\Pi \; p\text{exp} \rightarrow \text{pred}. \\
&(\Pi \; v : \text{exp}. \; \text{pf}(p \; v)) \rightarrow \text{pf } (\text{all } p)
\end{aligned}$$

## CERTIFICATE SIZE: EMPIRICAL DATA

One of the major problems with PCC is the size of the certificates.
Size of proof terms in Isabelle/HOL:

| Example | Size of proof term (lines) | Size of proof term (constructors) | Size of proof script (lines) |
|---------|------|------|------|
| AllImpl | 6 | 31 | 8 |
| AllExists | 6 | 26 | 7 |
| Arith | 295 | 1250 | 2 |

AllImpl :    $\forall AB.\ (A \wedge B \longrightarrow B \wedge A)$
AllExists :  $(\forall P.\ (\exists x.\ \forall y.\ P\ x\ y) \longrightarrow (\forall y.\ \exists x.\ P\ x\ y))$
Arith :      $\forall (m :: \mathrm{nat}).\ m < m + 1$

## CERTIFICATE SIZE: EMPIRICAL DATA

One of the major problems with PCC is the size of the certificates.
Size of proof terms in Isabelle/HOL:

| Example | Size of proof term (lines) | Size of proof term (constructors) | Size of proof script (lines) |
|---|---|---|---|
| const | 6 | 16 | 3 |
| cons with clarsimp | 31 | 136 | 3 |
| swap | 34819 | 137671 | 15 |
| count-down | 8584 | 25334 | 17 |
| list-reversal | 44082 | 162813 | 114 |

const : $\triangleright$ expr.Int $1$ : $\{(E, h, h', v, p).\ h' = h \wedge v = \text{IVal } 1 \wedge p = \langle (Suc\ 0)\ 0\ 0\ 0 \rangle$
swap : $\triangleright$ CALL swap : *spectable* swap
count : $\triangleright$ MH_InvokeStatic KountClass kount : *Mspectable* KountClass kount
rev : $\triangleright$ CALL rev : *spectable* rev

## DEEP VS SHALLOW EMBEDDING

When formalising a logic, how shall we represent assertions?

## Deep vs Shallow Embedding

When formalising a logic, how shall we represent assertions?

- **Deep Embedding**: define an explicit data structure of assertions

  ```
  data assn = true | false | and assn assn | ...
  ```

  Define an evaluation function that interprets an assertion
  *eval* : *state* ⇒ *assn* ⇒ *value*

# DEEP VS SHALLOW EMBEDDING

When formalising a logic, how shall we represent assertions?

- **Deep Embedding**: define an explicit data structure of assertions

  ```
  data assn = true | false | and assn assn | ...
  ```

  Define an evaluation function that interprets an assertion
  *eval* : *state* $\Rightarrow$ *assn* $\Rightarrow$ *value*

- **Shallow Embedding**: define assertions as functions over the state
  *type assn* $=$ *state* $\Rightarrow$ *value*

# DEEP VS SHALLOW EMBEDDING

When formalising a logic, how shall we represent assertions?

- **Deep Embedding**: define an explicit data structure of assertions

  ```
  data assn = true | false | and assn assn | ...
  ```

  Define an evaluation function that interprets an assertion
  $eval : state \Rightarrow assn \Rightarrow value$

- **Shallow Embedding**: define assertions as functions over the state
  $type\ assn = state \Rightarrow value$

Deep embeddings are usually easier to deal with.
Meta-properties over assertions may be harder to prove, though.

## STYLES OF PROGRAM LOGICS

Two styles of program logics have been proposed.

## STYLES OF PROGRAM LOGICS

Two styles of program logics have been proposed.

- Hoare-style logics: $\{P\}e\{Q\}$
  Assertions are parameterised over the "current" state.
  Example: Specification of an exponential function

  $$\{0 \le y \ \wedge \ x = X \wedge \ y = Y\} \ \exp(x, y) \ \{r = X^Y\}$$

  Note: $X, Y$ are **auxiliary variables** and must not appear in $e$

## STYLES OF PROGRAM LOGICS

Two styles of program logics have been proposed.

- Hoare-style logics: $\{P\}e\{Q\}$
  Assertions are parameterised over the "current" state.
  Example: Specification of an exponential function

$$\{0 \leq y \ \wedge \ x = X \wedge \ y = Y\} \exp(x, y) \ \{r = X^Y\}$$

  Note: $X, Y$ are **auxiliary variables** and must not appear in $e$

- VDM-style logics: $e : P$
  Assertions are parameterised over pre- and post-state.
  Because we have both pre- and post-state in the
  post-condition we do not need a separate pre-condition.
  Example: Specification of an exponential function

$$\{0 \leq y\} \exp(x, y) \ \{r = \grave{x}^{\grave{y}}\}$$

## A SIMPLE WHILE-LANGUAGE

Language:

$$
\begin{aligned}
e \quad ::= \quad & \text{skip} \\
& | \quad x := t \\
& | \quad e_1 ; e_2 \\
& | \quad \text{if } b \text{ then } e_1 \text{ else } e_2 \\
& | \quad \text{while } b \text{ do } e \\
& | \quad \text{call}
\end{aligned}
$$

# A Simple while-language

Language:

$$
\begin{aligned}
e \;::=\; & \text{skip} \\
& |\quad x := t \\
& |\quad e_1 ; e_2 \\
& |\quad \text{if } b \text{ then } e_1 \text{ else } e_2 \\
& |\quad \text{while } b \text{ do } e \\
& |\quad \text{call}
\end{aligned}
$$

A judgement has this form (for now!)

$$\vdash \{P\}\ e\ \{Q\}$$

A judgement is valid if the following holds

$$\forall z\ s\ t.\ s \overset{e}{\leadsto} t \Rightarrow P\ z\ s \Rightarrow Q\ z\ t$$

# A SIMPLE HOARE-STYLE LOGIC

$$\overline{\vdash \{P\} \; \texttt{skip} \; \{P\}} \quad \text{(SKIP)} \qquad \overline{\vdash \{\lambda z \; s. \; P \; z \; s[t/x]\} \; x := t \; \{P\}}$$
$$\text{(ASSIGN)}$$

$$\frac{\vdash \{P\} \; e_1 \; \{R\} \quad \{R\} \; e_2 \; \{Q\}}{\vdash \{P\} \; e_1;e_2 \; \{Q\}} \qquad \text{(COMP)}$$

$$\frac{\vdash \{\lambda z \; s. \; P \; z \; s \; \wedge \; b \; s\} \; e_1 \; \{Q\} \quad \vdash \{\lambda z \; s. \; P \; z \; s \; \wedge \; \neg(b \; s)\} \; e_2 \; \{Q\}}{\vdash \{P\} \; \texttt{if} \; b \; \texttt{then} \; e_1 \; \texttt{else} \; e_2 \{Q\}} \; \text{(IF)}$$

$$\frac{\vdash \{\lambda z \; s. \; P \; z \; s \; \wedge \; b \; s\} \; e \; \{P\}}{\vdash \{P\} \; \texttt{while} \; b \; \texttt{do} \; e\{\lambda z \; s. \; P \; z \; s \; \wedge \; \neg(b \; s)\}} \qquad \text{(WHILE)}$$

$$\frac{\vdash \{P\} \; body \; \{Q\}}{\vdash \{P\} \; \texttt{CALL} \; \{Q\}} \qquad \text{(CALL)}$$

# A SIMPLE HOARE-STYLE LOGIC (STRUCTURAL RULES)

The consequence rule allows us to weaken the pre-condition and to strengthen the post-condition:

$$\frac{\forall s\ t.\ (\forall z.\ P'\ z\ s \Rightarrow P\ z\ s) \quad \vdash \{P'\}\ e\ \{Q'\} \quad \forall s\ t.\ (\forall z.\ Q\ z\ s \Rightarrow Q'\ z\ s)}{\vdash \{P\}\ e\ \{Q\}}$$

$$(\text{CONSEQ})$$

## RECURSIVE FUNCTIONS

In order to deal with recursive functions, we need to collect the knowledge about the behaviour of the functions.

We extend the judgement with a context $\Gamma$, mapping expressions to Hoare-Triples:

$$\Gamma \vdash \{P\} \ e \ \{Q\}$$

where $\Gamma$ has the form $\{\ldots, (P', e', Q'), \ldots\}$.

## RECURSIVE FUNCTIONS

Now, the call rule for recursive, parameter-less functions looks like this:

$$\frac{\Gamma \cup \{(P, \texttt{CALL}, Q)\} \vdash \{P\} \ body \ \{Q\}}{\Gamma \vdash \{P\} \ \texttt{CALL} \ \{Q\}} \qquad (\text{CALL})$$

We collect the knowledge about the (one) function in the context, and prove the body.

**Note**: This is a rule for partial correctness: for total correctness we need some form of measure.

## Recursive Functions

To extract information out of the context we need and axiom rule

$$\frac{(P, e, Q) \in \Gamma}{\Gamma \vdash \{P\} \ e \ \{Q\}} \qquad \text{(ax)}$$

## RECURSIVE FUNCTIONS

To extract information out of the context we need and axiom rule

$$\frac{(P, e, Q) \in \Gamma}{\Gamma \vdash \{P\} \ e \ \{Q\}} \quad \text{(AX)}$$

Note that we now use a **Gentzen-style** logic (one with contexts) rather than a Hilbert-style logic.

## MORE TROUBLES WITH RECURSIVE FUNCTIONS

Assume we have this simple recursive program:

```
if i=0 then skip else i := i-1 ; call ; i := i+1
```

## MORE TROUBLES WITH RECURSIVE FUNCTIONS

Assume we have this simple recursive program:

```
if i=0 then skip else i := i-1 ; call ; i := i+1
```

The proof of $\{i = N\}$ call $\{i = N\}$ proceeds as follows

$$\frac{\rule{4cm}{0.4pt}}{\vdash \{i = N\} \text{ CALL } \{i = N\}}$$

## MORE TROUBLES WITH RECURSIVE FUNCTIONS

Assume we have this simple recursive program:

```
if i=0 then skip else i := i-1 ; call ; i := i+1
```

The proof of $\{i = N\}$ call $\{i = N\}$ proceeds as follows

$$\frac{\{(i = N, \text{CALL}, i = N)\} \vdash \{i = N\} \ \text{i} := \text{i} - 1; \text{CALL}; \text{i} := \text{i} + 1 \ \{i = N\}}{\vdash \{i = N\} \ \text{CALL} \ \{i = N\}}$$

## MORE TROUBLES WITH RECURSIVE FUNCTIONS

Assume we have this simple recursive program:

```
if i=0 then skip else i := i-1 ; call ; i := i+1
```

The proof of $\{i = N\}$ call $\{i = N\}$ proceeds as follows

$$\frac{\dfrac{\{(i = N, \text{CALL}, i = N)\} \vdash \; \{i = N - 1\} \; \text{CALL} \; \{i = N - 1\}}{\{(i = N, \text{CALL}, i = N)\} \vdash \; \{i = N\} \; \texttt{i} := \texttt{i} - 1; \text{CALL}; \texttt{i} := \texttt{i} + 1 \; \{i = N\}}}{\vdash \; \{i = N\} \; \text{CALL} \; \{i = N\}}$$

## MORE TROUBLES WITH RECURSIVE FUNCTIONS

Assume we have this simple recursive program:

```
if i=0 then skip else i := i-1 ; call ; i := i+1
```

The proof of $\{i = N\}$ `call` $\{i = N\}$ proceeds as follows

$$\frac{\{(i = N, \text{CALL}, i = N)\} \vdash \{i = N - 1\} \text{ CALL } \{i = N - 1\}}{\dfrac{\{(i = N, \text{CALL}, i = N)\} \vdash \{i = N\} \ \text{i} := \text{i} - 1; \text{CALL}; \text{i} := \text{i} + 1 \ \{i = N\}}{\vdash \{i = N\} \text{ CALL } \{i = N\}}}$$

But how can we prove $\{i = N - 1\}\text{CALL}\{i = N - 1\}$ from
$\{i = N\}\text{CALL}\{i = N\}$?

## More Troubles with Recursive Functions

Assume we have this simple recursive program:

```
if i=0 then skip else i := i-1 ; call ; i := i+1
```

The proof of $\{i = N\}$ call $\{i = N\}$ proceeds as follows

$$\frac{\dfrac{\{(i = N, \text{CALL}, i = N)\} \vdash \{i = N - 1\} \text{ CALL } \{i = N - 1\}}{\{(i = N, \text{CALL}, i = N)\} \vdash \{i = N\} \text{ i} := \text{i} - 1; \text{CALL}; \text{i} := \text{i} + 1 \{i = N\}}}{\vdash \{i = N\} \text{ CALL } \{i = N\}}$$

But how can we prove $\{i = N - 1\}$CALL$\{i = N - 1\}$ from
$\{i = N\}$CALL$\{i = N\}$?
We need to **instantiate** $N$ with $N - 1$!

## RECURSIVE FUNCTIONS

To be able to instantiate auxiliary variables we need a more powerful consequence rule:

$$\frac{\Gamma \vdash \{P'\}\ e\ \{Q'\} \quad \forall s\ t.\ (\forall z.\ P'\ z\ s \Rightarrow Q'\ z\ t)\ \Rightarrow\ (\forall z.\ P\ z\ s \Rightarrow Q\ z\ t)}{\Gamma \vdash \{P\}\ e\ \{Q\}}$$

$$(\text{CONSEQ})$$

Now we are allowed to proof $P \Rightarrow Q$ under the knowledge that we can choose $z$ freely as long as $P' \Rightarrow Q'$ is true.
This complex rule for **adaptation** is one of the main disadvantages of Hoare-style logics.

## Extending the Logic with Termination

The Call and While rules need to use a well-founded ordering $<$ and a side condition saying that the body is smaller w.r.t. this ordering:

$$
\begin{array}{c}
wf\ < \\
\forall s'.\ \{(\lambda z\ s.P\ z\ s\ \wedge\ \ s < s', \mathtt{CALL}, Q)\} \\
\vdash_T \{\lambda z\ s.P\ z\ s\ \wedge\ \ s = s'\} body\ \{Q\} \\
\hline
\vdash_T \{P\}\ \mathtt{CALL}\{Q\}
\end{array}
$$

Note the explicit quantification over the state s'. Read it like this

*The pre-state s must be smaller than a state $s'$, which is the post-state.*

## EXTENDING THE LOGIC WITH MUTUAL RECURSION

To cover mutual recursion a different derivation system $\vdash_M$ is defined.
Judgements in $\vdash_M$ are extended to sets of Hoare triples, informally:

$$\Gamma \vdash_M \{(P_1, e_1, Q_1), \ldots, (P_n, e_n, Q_n)\}$$

The Call rule is generalised as follows

$$\frac{\bigcup p. \{(P\ p, \text{CALL p}, Q\ p)\} \vdash_M \bigcup p. \{(P\ p, body\ p, Q\ p)\}}{\emptyset \vdash_M \bigcup p. \{(P\ p, \text{CALL p}, Q\ p)\}}$$

# FURTHER READING

📕 Thomas Kleymann, *Hoare Logic and VDM: Machine-Checked Soundness and Completeness Proofs*, Lab. for Foundations of Computer Science, Univ of Edinburgh, LFCS report ECS-LFCS-98-392, 1999.
http://www.lfcs.informatics.ed.ac.uk/reports/98/ECS-LFCS-98-3

📕 Tobias Nipkow, *Hoare Logics for Recursive Procedures and Unbounded Nondeterminism*, in CSL 2002 — Computer Science Logic, LNCS 2471, pp. 103–119, Springer, 2002.

## CHALLENGE: MINIMISING THE TCB

This aspect is the emphasis of the **Foundational PCC** approach.

An infrastructure developed by the group of Andrew Appel at Princeton [1].

**Motivation**: With complex logics and VCGs, there is a big danger of introducing bugs in software that needs to be trusted.

# The Philosophy of Foundational PCC

Define safety policy directly on the **operational semantics** of the code.

Certificates are proofs over the operational semantics.

It minimises the TCB because no trusted verification condition generator is needed.

Pros and cons:

- ☺ **more flexible**: not restricted to a particular type system as the language in which the proofs are phrased;
- ☺ **more secure**: no reliance on VCG.
- ☹ **larger proofs**

## Conventional vs Foundational PCC

Re-examine the logic for memory safety, eg.

$$\frac{m \vdash e : \tau \; list \quad e \neq 0}{\begin{array}{l} m \vdash e : addr \; \wedge \; m \vdash e + 4 : addr \; \wedge \\ m \vdash sel(m, e) : \tau \; \wedge \; m \vdash sel(m, e + 4) : \tau \; list \end{array}}$$
(LISTELIM)

The rule has **built-in knowledge about the type-system**, in this case representing the data layout of the compiler ("*Type specialised PCC*") $\implies$ dangerous if soundness of the logic is not checked mechanically!

## Logic rules in Foundational PCC

In foundational PCC the rules work on the operational semantics:

$$\frac{m \models e : \tau\ list \quad e \neq 0}{\begin{array}{l} m \models e : addr\ \wedge\ m \models e + 4 : addr\ \wedge \\ m \models sel(m, e) : \tau\ \wedge\ m \models sel(m, e + 4) : \tau\ list \end{array}} \quad (\textsc{ListElim})$$

This looks similar to the previous rule but has a very different meaning: $\models$ is a predicate over the formal model of the computation, and the above rule can be proven as a lemma, $\vdash$ is an encoding of a type-system on top of the operational semantics and thus needs a **soundness proof**.

# COMPONENTS OF A FOUNDATIONAL PCC INFRASTRUCTURE

Operational semantics and safety properties are directly encoded in a **higher-order logic**.

As language for the certificates, the LF metalogic framework is used.

For development and for proof-checking the Twelf theorem proofer is used.

## SPECIFYING SAFETY

To specify safety, the operational semantics is written in such a way, that it gets stuck whenever the safety condition is violated.

## SPECIFYING SAFETY

To specify safety, the operational semantics is written in such a way, that it gets stuck whenever the safety condition is violated.

Example: operational semantics on assembler code.
Safety policy: "only readable addresses are loaded".
Define a predicate: $readable(x) \equiv 0 \leq x \leq 1000$

## SPECIFYING SAFETY

To specify safety, the operational semantics is written in such a
way, that it gets stuck whenever the safety condition is violated.

Example: operational semantics on assembler code.
Safety policy: "only readable addresses are loaded".
Define a predicate: $readable(x) \equiv 0 \leq x \leq 1000$
The semantics of a load operation LD $r_i, c(r_j)$ is now written as
follows:

$$
\begin{aligned}
load(i, j, c) \quad \equiv \quad & \lambda \ r \ m \ r' \ m'. \\
& r'(i) = m(r(j) + c) \ \wedge \ readable(r(j) + c) \ \wedge \\
& (\forall x \neq i. \ r'(x) = r(x)) \ \wedge \ m' = m
\end{aligned}
$$

## SPECIFYING SAFETY

To specify safety, the operational semantics is written in such a way, that it gets stuck whenever the safety condition is violated.

Example: operational semantics on assembler code.
Safety policy: "only readable addresses are loaded".
Define a predicate: $readable(x) \equiv 0 \leq x \leq 1000$
The semantics of a load operation LD $r_i, c(r_j)$ is now written as follows:

$$
\begin{aligned}
load(i, j, c) \quad \equiv \quad & \lambda \ r \ m \ r' \ m'. \\
& r'(i) = m(r(j) + c) \ \wedge \ readable(r(j) + c) \ \wedge \\
& (\forall x \neq i. \ r'(x) = r(x)) \ \wedge \ m' = m
\end{aligned}
$$

**Note:** the clause for nothing else changes, quickly becomes awkward when doing these proofs
$\implies$ Separation Logic (Reynolds'02) tackles this problem.

## Main issues in FPCC

The main task in this framework becomes the **semantic modelling of types**: indexed semantic model to describe contravariant types, eg. $e = \text{APP } of\ e\ e \mid \text{LAM } of\ e \rightarrow e$

Naive model: type $=$ set of values

Indexed model: type $=$ set of $< k, v >$, where $k$ is an approximation index, $v$ is a value
$< k, v > \in \tau$ means $v$ has approximate type $\tau$ and programs running less than $k$ steps can't tell a difference
$\implies$ induction principle over steps of execution

# FURTHER READING

📕 Andrew Appel, *Foundational Proof-Carrying Code* in LICS'01
  — Symposium on Logic in Computer Science, 2001.
  http://www.cs.princeton.edu/~appel/papers/fpcc.pdf

## CCURED

A system for checking **pointer-safety** of C programs, developed by the group of George Necula at Berkeley.

Uses a hybrid mechanism of static type checking and run-time checks.

**Goal:** Prove pointer safety statically, where possible, and minimise required run-time checks.

## THE CCURED TYPE SYSTEM

Extension of the standard C type system with extension for
**pointers into arrays and dynamic types**.

Efficient type inference is possible and demonstrated for this type
system.

## THE CORE LANGUAGE

Mini-C language:

$$e ::= x \mid n \mid e_1 \text{ op } e_2 \mid (\tau)e \mid e_1 \oplus e_2 \mid !e$$
$$c ::= \text{skip} \mid c_1;c_2 \mid e_1 := e_2$$

# THE CCURED TYPE SYSTEM: POINTERS

C contains 2 evil pointer operations: arithmetic and casts.

The type system distinguishes between 3 kinds of pointers:

- **Safe pointers**: no arithmetic or casts; represented as an address
- **Sequence pointers**: arithmetic but **no casts**; represented as a region
- **Dynamic pointers**: **casts**, all bets are off! represented as a region

## EXAMPLE PROGRAM

Sum over an array of boxed integers:

```
int **a;  /* array */         int i;      // index
int acc;  /* accumulator */   int **p;    // elem ptr
int *e;   /* unboxer */
acc = 0;
for (i=0; i<100; i++) {
  p = a + i;                // ptr arithm
  e = *p;                   // read elem
  while ((int)e % 2 == 0) { // check tag
    e = *(int **)e;         // unbox
  }
  acc += ((int)e >> 1);     // strip tag
}
```

## EXAMPLE PROGRAM

Sum over an array of boxed integers:

```
int **a;  /* array */       int i;      // index
int acc; /* accumulator */  int **p;    // elem ptr
int *e;   /* unboxer */
acc = 0;
for (i=0; i<100; i++) {
  p = a + i;                // ptr arithm
  e = *p;                   // read elem
  while ((int)e % 2 == 0) { // check tag
    e = *(int **)e;         // unbox
  }
  acc += ((int)e >> 1);     // strip tag
}
```

a and p point into an array with elems of type int *

## EXAMPLE PROGRAM

Sum over an array of boxed integers:

```
int **a;  /* array */        int i;      // index
int acc;  /* accumulator */  int **p;    // elem ptr
int *e;   /* unboxer */
acc = 0;
for (i=0; i<100; i++) {
  p = a + i;                 // ptr arithm
  e = *p;                    // read elem
  while ((int)e % 2 == 0) {  // check tag
    e = *(int **)e;          // unbox
  }
  acc += ((int)e >> 1);      // strip tag
}
```

a is subject to pointer arithm ("sequence pointer")
$\Longrightarrow$ check for out of bounds

## EXAMPLE PROGRAM

Sum over an array of boxed integers:

```
int **a;  /* array */        int i;       // index
int acc; /* accumulator */  int **p;     // elem ptr
int *e;   /* unboxer */
acc = 0;
for (i=0; i<100; i++) {
  p = a + i;               // ptr arithm
  e = *p;                  // read elem
  while ((int)e % 2 == 0) { // check tag
    e = *(int **)e;         // unbox
  }
  acc += ((int)e >> 1);    // strip tag
}
```

p has no arithmetic ("safe pointer")
$\implies$ no bounds check needed

## EXAMPLE PROGRAM

Sum over an array of boxed integers:

```
int **a;  /* array */       int i;       // index
int acc; /* accumulator */  int **p;     // elem ptr
int *e;   /* unboxer */
acc = 0;
for (i=0; i<100; i++) {
  p = a + i;                // ptr arithm
  e = *p;                   // read elem
  while ((int)e % 2 == 0) { // check tag
    e = *(int **)e;         // unbox
  }
  acc += ((int)e >> 1);     // strip tag
}
```

e is subject to a type cast ("dynamic pointer")
$\implies$ nothing known about underlying type

## SAFE POINTERS

Invariant for SAFE pointers:

> A **SAFE** pointer to type T is either 0 or else it points to
> a valid area of memory containing an object of type T.
> Furthermore, all other pointers to the same area are also
> SAFE and agree on the type T of the stored object.

Run-time check: null-pointer reference.

## SEQUENCE POINTERS

Invariants for Sequence pointers:

- Cannot be cast (passing actual arguments and returning are implicit casts).

- Can be subject to pointer arithmetic (adding or subtracting an integer from it).

- Can be set to any integer value.

- Can be cast to an integer and can be subtracted from another pointer (useful for comparisons).

- Sequence pointers are represented using three words.

Run-time checks: null-pointer check and bounds check.

## OPERATIONAL SEMANTICS

The value of an integer, or a safe pointer is an integer $n$; the value of a sequence or dynamic pointer is a **home**, modelled as a pair $\mathbb{N} \times \mathbb{N}$ of start address and offset.

$$v ::= n \mid \langle h, n \rangle$$

## OPERATIONAL SEMANTICS

The value of an integer, or a safe pointer is an integer $n$; the value of a sequence or dynamic pointer is a **home**, modelled as a pair $\mathbb{N} \times \mathbb{N}$ of start address and offset.

$$v ::= n \mid \langle h, n \rangle$$

Each home is tagged as being an integer or a pointer, and has an associated **kind** and **size** functions. The semantic domain for pointers:

$$
\begin{aligned}
\| \text{ int } \|_H &= \mathbb{N} \\
\| \text{ DYNAMIC } \|_H &= \{\langle h, n \rangle \mid h \in H \wedge (h = 0 \vee \text{kind}(h) = untyped)\} \\
\| \tau \text{ ref SEQ } \|_H &= \{\langle h, n \rangle \mid h \in H \wedge (h = 0 \vee \text{kind}(h) = typed(\tau))\} \\
\| \tau \text{ ref SAFE } \|_H &= \{h + i \mid h \in H \wedge 0 \leq i \leq size(h) \wedge \\
&\qquad (h = 0 \vee \text{kind}(h) = typed(\tau))\}
\end{aligned}
$$

# OPERATIONAL SEMANTICS (POINTERS)

$$\frac{\Sigma, M \vdash e_1 \Downarrow \langle h, n_1 \rangle \quad \Sigma, M \vdash e_2 \Downarrow n_2}{\Sigma, M \vdash e_1 \oplus e_2 \Downarrow \langle h_1, n_1 + n_2 \rangle}$$
$$(\text{POINTER ARTIHM})$$

$$\frac{\Sigma, M \vdash e \Downarrow \langle h, n \rangle}{\Sigma, M \vdash (\texttt{int})e \Downarrow h + n} \qquad (\text{CASTTOINT})$$

$$\frac{\Sigma, M \vdash e \Downarrow n}{\Sigma, M \vdash (\tau \texttt{ ref SEQ})e \Downarrow \langle 0, n \rangle} \qquad (\text{CASTTOSEQ})$$

$$\frac{\Sigma, M \vdash e \Downarrow \langle h, n \rangle \quad \mathbf{0 \leq n \leq \texttt{size}(\mathbf{h})}}{\Sigma, M \vdash (\tau \texttt{ ref SAFE})e \Downarrow h + n}$$
$$(\text{CASTTOSAFE})$$

## Operational Semantics (read operations)

Two kinds of reads, with different obligations for run-time checks:

$$\frac{\Sigma, M \vdash e \Downarrow n \quad \mathbf{n \neq 0}}{\Sigma, M \vdash !e \Downarrow M(n)} \tag{SafeRd}$$

$$\frac{\Sigma, M \vdash e \Downarrow \langle h, n \rangle \quad \mathbf{h \neq 0} \quad 0 \leq n \leq \texttt{size}(h)}{\Sigma, M \vdash !e \Downarrow M(h + n)} \tag{DynRd}$$

$$\frac{\Sigma, M \vdash e_1 \Downarrow n \quad \mathbf{n \neq 0} \quad \Sigma, M \vdash e_2 \Downarrow v}{\Sigma, M \vdash e_1 := e_2 \Downarrow M(n \mapsto v)} \tag{SafeWr}$$

$$\frac{\Sigma, M \vdash e_1 \Downarrow \langle h, n \rangle \quad \mathbf{h \neq 0} \quad 0 \leq n \leq \texttt{size}(h) \quad \Sigma, M \vdash e_2 \Downarrow v}{\Sigma, M \vdash e_1 := e_2 \Downarrow M(h + n \mapsto v)} \tag{DynWr}$$

Hans-Wolfgang Loidl          Proof-Carrying-Code

## THE CCURED TYPE SYSTEM: RULES

The type system keeps track of the kind of pointers.
Rules for converting pointers:

$$\frac{}{\tau \leq \tau} \qquad\qquad \frac{}{\tau \leq \mathtt{int}} \qquad\qquad \frac{}{\mathtt{int} \leq \tau \; \mathtt{ref} \; \mathtt{SEQ}}$$

$$\frac{}{\mathtt{int} \leq \mathtt{DYNAMIC}}$$

$$\frac{}{\tau \; \mathtt{ref} \; \mathtt{SEQ} \leq \tau \; \mathtt{ref} \; \mathtt{SAFE}}$$

## TYPING RULES FOR COMMANDS

$$\frac{}{\Gamma \vdash \mathtt{skip}} \qquad \frac{\Gamma \vdash c_1 \quad \Gamma \vdash c_2}{\Gamma \vdash c_1; c_2} \qquad \frac{\Gamma \vdash e : \tau \ \mathtt{ref}\ \mathtt{SAFE} \quad \Gamma \vdash e' : \tau}{\Gamma \vdash e := e'}$$

$$\frac{\Gamma \vdash e : \mathtt{DYNAMIC} \quad \Gamma \vdash e' : \mathtt{DYNAMIC}}{\Gamma \vdash e := e'}$$

## TYPING RULES FOR EXPRESSIONS

$$\frac{\Gamma(x) = \tau}{\Gamma \vdash x : \tau} \qquad \frac{\Gamma \vdash e_1 : \texttt{int} \quad \Gamma \vdash e_2 : \texttt{int}}{\Gamma \vdash e_1 \text{ op } e_2 : \texttt{int}} \qquad \frac{\Gamma \vdash e : \tau' \quad \tau' \leq \tau}{\Gamma \vdash (\tau)e : \tau}$$

$$\frac{}{\Gamma \vdash (\tau \text{ ref SAFE})0 : \tau \text{ ref SAFE}} \qquad \frac{\Gamma \vdash e_1 : \tau \text{ ref SEQ} \quad \Gamma \vdash e_2 : \texttt{int}}{\Gamma \vdash e_1 \oplus e_2 : \tau \text{ ref SEQ}}$$

$$\frac{\Gamma \vdash e_1 : \texttt{DYNAMIC} \quad \Gamma \vdash e_2 : \texttt{int}}{\Gamma \vdash e_1 \oplus e_2 : \texttt{DYNAMIC}} \qquad \frac{\Gamma \vdash e : \tau \text{ ref SAFE}}{\Gamma \vdash !e : \tau} \qquad \frac{\Gamma \vdash e : \texttt{DYNAMIC}}{\Gamma \vdash !e : \texttt{DYNAMIC}}$$

$\Sigma, M_H \vdash e \Downarrow \textit{CheckFailed}$ means a run-time check failed during the execution of expression $e$.

## THEOREMS

$\Sigma, M_H \vdash e \Downarrow CheckFailed$ means a run-time check failed during the execution of expression $e$.

### THEOREM (PROGRESS AND TYPE PRESERVATION)

If $\Gamma \vdash e : \tau$ and $\Sigma \in \| \Gamma \|_H$ and $M$ is well-formed, then either $\Sigma, M_H \vdash e \Downarrow CheckFailed$ or $\Sigma, M_H \vdash e \Downarrow v$ and $v \in \| \tau \|_H$.

## THEOREMS

$\Sigma, M_H \vdash c \implies CheckFailed$ means a run-time check failed during the execution of command $c$.

## THEOREMS

$\Sigma, M_H \vdash c \implies$ *CheckFailed* means a run-time check failed during the execution of command $c$.

### THEOREM (PROGRESS FOR COMMANDS)

*If $\Gamma \vdash c$ and $\Sigma \in \| \Gamma \|_h$ and $M_H$ is well-formed then either $\Sigma, M_H \vdash c \implies$ CheckFailed or $\Sigma, M_H \vdash c \implies M'_H$ and $M'_H$ is well-formed.*

# Main results

- An efficient inference algorithm attaches
  ref SEQ, ref SAFE, DYNAMIC annotations to plain C code.
- Most of the checks can be done statically.
- The performance overhead of the remaining run-time checks
  is moderate: 0–150%

# FURTHER READING

📕 *CCured: Type-Safe Retrofitting of Legacy Code*, in POPL'02 — ACM Symposium on Principles of Programming Languages, 2002.